

Multi-year Predictions of North Atlantic Hurricane Frequency: Promise and limitations

Gabriel A. Vecchi¹, Rym Msadek¹, Whit Anderson¹, You-Soon Chang¹, Thomas Delworth¹, Keith Dixon¹, Rich Gudgel¹, Anthony Rosati¹, Bill Stern¹, Gabriele Villarini², Andrew Wittenberg¹, Xiasong Yang¹, Fanrong Zeng¹, Rong Zhang¹, Shaoqing Zhang¹

1. Geophysical Fluid Dynamics Laboratory, NOAA, Princeton, NJ, USA
2. IIHR-Hydroscience & Engineering, The University of Iowa, Iowa City, IA, USA

For submission to *J. Climate*

**MANUSCRIPT IN REVIEW,
DO NOT QUOTE OR CITE WIHTOUT PERMISSION.**

July 19, 2012

Corresponding Author:
Gabriel A. Vecchi
Geophysical Fluid Dynamics Laboratory / NOAA
US Route 1, Forrester Campus
Princeton, NJ 08542
Tel: (609) 452-6583, Fax: (609) 987-5063, email: gabriel.a.vecchi@noaa.gov

1 **Abstract**

2 Retrospective predictions of multi-year North Atlantic hurricane frequency are explored,
3 by applying a hybrid statistical-dynamical forecast system to initialized and non-
4 initialized multi-year forecasts of tropical Atlantic and tropical mean sea surface
5 temperatures (SSTs) from two global climate model forecast systems. By accounting for
6 impacts of initialization and radiative forcing, retrospective predictions of five-year mean
7 and nine-year mean tropical Atlantic hurricane frequency show significant correlation
8 relative to a null hypothesis of zero correlation. The retrospective correlations are
9 increased in a two-model average forecast and by using a lagged-ensemble approach,
10 with the two-model ensemble decadal forecasts hurricane frequency over 1961-2011
11 yielding correlation coefficients that approach 0.9.

12 These encouraging retrospective multi-year hurricane predictions, however, should be
13 interpreted with care: although initialized forecasts have higher nominal skill than
14 uninitialized ones, the relatively short record and large autocorrelation of the time series
15 limits our confidence in distinguishing between the skill due to external forcing and that
16 added by initialization. The nominal increase in correlation in the initialized forecasts
17 relative to the uninitialized experiments is due to improved representation of the multi-
18 year tropical Atlantic SST anomalies. The skill in the initialized forecasts comes in large
19 part from the persistence of a mid-1990s shift by the initialized forecasts, rather than
20 from predicting its evolution. Predicting shifts like that observed in 1994-1995 remains a
21 critical issue for the success of multi-year forecasts of Atlantic hurricane frequency. The
22 retrospective forecasts highlight the possibility that changes in observing system impact
23 forecast performance.

1

2 **I. Introduction**

3 Predicting and projecting future North Atlantic hurricane activity is a topic of
4 scientific interest (*e.g.*, Gray 1984; Knutson and Tuleya 2004; Emanuel 2005; Camargo
5 *et al.* 2007a; Vecchi *et al.* 2008; Smith *et al.* 2010; Knutson *et al.* 2010; Vecchi *et al.*
6 2011; Villarini *et al.* 2011.a; Villarini and Vecchi 2012b-d) and high societal significance
7 (Pielke Jr. *et al.* 2008; Mendelsohn *et al.* 2012; Peduzzi *et al.* 2012). Seasonal basin-wide
8 frequency of North Atlantic hurricanes has exhibited variability on a variety of
9 timescales, from interannual to multi-decadal, although it remains unclear whether there
10 has been any century-scale trend in Atlantic hurricane frequency (*e.g.*, Mann and
11 Emanuel 2006; Vecchi and Knutson 2008, 2011; Landsea *et al.* 2011; Villarini *et al.*
12 2011b).

13 The scientific basis for predictions of seasonal hurricane activity at leads of one to
14 three seasons has been developed (*e.g.*, Gray 1984; Elsner and Jagger 2006; Vitart 2006;
15 Camargo *et al.* 2007a,b; Vitart *et al.* 2007; Klotzbach and Gray 2009; Wang *et al.* 2009;
16 Kim and Webster 2010; LaRow *et al.* 2010; Zhao *et al.* 2010; Alessandri *et al.* 2011;
17 Chen and Lin 2011; Vecchi *et al.* 2011; Villarini and Vecchi 2012d), leading to the
18 identification of different potential sources of skill, both local and remote.

19 Decadal to centennial projections of seasonal hurricane activity in response to
20 changes in external forcing (greenhouse gases, aerosols, volcanoes, and solar) have been
21 made (*e.g.* Oouchi *et al.* 2006; Knutson *et al.* 2008; Emanuel *et al.* 2008; Gualdi *et al.*
22 2008; Vecchi *et al.* 2008; Sugi *et al.* 2009, 2012; Zhao *et al.* 2009; Bender *et al.* 2010;
23 Knutson *et al.* 2010; Knutson *et al.* 2010; Villarini *et al.* 2011a; Zhao and Held 2011;

Villarini and Vecchi 2012b,c). The basis for these projections is the possibility that radiatively-forced climate change could influence the climatic conditions to which hurricanes are sensitive, such as large-scale circulation, wind shear, ocean temperatures, potential intensity and humidity (*e.g.*, Emanuel 1987, 2007; Broccoli and Manabe 1990; Shen *et al.* 2000; Knutson and Tuleya 2004; Camargo *et al.* 2007b; Vecchi and Soden 2007a,b). Recent model results span a relatively wide range of possibilities for North Atlantic hurricane frequency (including increases or decreases) under enhanced CO₂-induced warming, while there is a wider tendency for hurricane intensity to increase in these studies (*e.g.*, Knutson and Tuleya 2004; Knutson *et al.* 2008, Emanuel *et al.* 2008, Gualdi *et al.* 2008; Knutson *et al.* 2008; Vecchi *et al.* 2008; Sugi *et al.* 2009, 2012; Zhao *et al.* 2009, Bender *et al.* 2010; Knutson *et al.* 2010, Villarini *et al.* 2011a; Villarini and Vecchi 2012b,c). There are indications that changes in atmospheric aerosols could influence past and projected hurricane activity, with increases (decreases) in Atlantic aerosol loading driving decreases (increases) in Atlantic hurricane activity (Mann and Emanuel 2006; Evan *et al.* 2009, Villarini and Vecchi 2012b,c).

Assessing hurricane predictability at intermediate timescales, between seasonal predictions and multi-decadal projections, is an emerging field of research. In addition to potential influences from changes in radiative forcing, internal variations of the climate system could play a large role in changes of hurricane frequency on timescales of decades (*e.g.*, Goldenberg *et al.* 1996; Zhang and Delworth 2006, 2009; Knight *et al.* 2006; Latif *et al.* 2007; Dunstone *et al.* 2011; Villarini *et al.* 2011; Villarini and Vecchi 2012b). There are physical reasons to expect coherent multi-year hurricane variations to be tied to ocean changes (*e.g.*, Goldenberg *et al.* 1996, Zhang and Delworth 2005, 2006, 2009;

1 Knight *et al.* 2006, Latif *et al.* 2007, Dunstone *et al.* 2011). There is also indication that
2 some of the relevant ocean changes may be potentially predictable on decadal timescales
3 (*e.g.*, Griffies and Bryan 1997a,b; Pohlmann *et al.* 2004; Collins *et al.* 2006; Msadek *et*
4 *al.* 2010; Teng *et al.* 2011; Yang *et al.* 2012; Rosati *et al.* 2012). As decadal variability
5 and the associated predictability can result from both internal and externally forced
6 fluctuations (*e.g.*, Rotstayn and Lohmann 2002; Hawkins and Sutton 2009; Chang *et al.*
7 2011a; Villarini *et al.* 2011; Booth *et al.* 2012; Villarini and Vecchi 2012b), one has to
8 consider skill arising from both external factors and internal variability on multi-year
9 timescales. A number of modeling groups are now following the same framework for the
10 Fifth Coupled Model Intercomparison Project (CMIP5; Taylor *et al.* 2012) to be assessed
11 as part of the 5th Assessment Report of the Intergovernmental Panel on Climate Change
12 (IPCC-AR5), by performing decadal predictions initialized with estimates of the
13 observed state of the climate system (Taylor *et al.* 2012, Meehl *et al.* 2012). While for sea
14 surface temperature (SST), most of the skill on multi-year timescales arises from
15 predicting the warming trend associated with radiative forcing changes (*e.g.*, van
16 Oldenborgh *et al.* 2012; Rosati *et al.* 2012), there is at least one study suggesting that
17 initialization can increase the skill in multi-year hurricane forecasts (Smith *et al.* 2010).
18 In this paper we explore the ability of a hybrid statistical-dynamical hurricane forecasting
19 system to retrospectively predict multi-year hurricane activity in the Atlantic using two
20 different coupled climate models, including the one used by Smith *et al.* (2010). We
21 explore the skill of North Atlantic hurricane frequency resulting from changing radiative
22 forcing and from natural variability. We assess the improvement in skill due to

1 initialization and discuss the source of this improved skill and its implications for future
2 multi-year forecasts of North Atlantic hurricane frequency.

3 **II. Data and Methods**

4 *A. Statistical hurricane emulator:*

5 We use a hybrid statistical-dynamical North Atlantic hurricane frequency prediction
6 framework to explore the predictability of multi-year hurricane activity. This framework
7 has been shown to exhibit retrospective skill in seasonal hurricane forecasts from as early
8 as boreal winter prior to the hurricane season (Vecchi *et al.* 2011). It combines a
9 statistical emulator of a high-resolution dynamical atmospheric model (Zhao *et al.* 2009,
10 2010) and initialized forecasts of SST. The statistical emulator is formulated as a Poisson
11 regression model with two predictors: Tropical Atlantic SST and Tropical-mean SST,
12 each averaged over the August-October season. The choice of these two predictors is
13 motivated by dynamical considerations, observed relationships between hurricane
14 activity and SST, and the sensitivity of dynamical models to SST perturbations (*e.g.*,
15 Shen *et al.* 2000; Sobel *et al.* 2002; Tang and Neelin 2004; Latif *et al.* 2007; Vecchi and
16 Soden 2007a; Gualdi *et al.* 2008; Knutson *et al.* 2008; Swanson 2008; Vecchi *et al.* 2008;
17 Ramsay and Sobel 2011; Villarini *et al.* 2010, 2011a, 2012; Camargo *et al.* 2012;
18 Villarini and Vecchi 2012a,d). Following Vecchi *et al.* (2011), we model the rate of
19 occurrence (λ ; the expected value of the aggregate seasonal number) of North Atlantic
20 hurricane frequency using a Poisson regression model as follows:

$$21 \quad \lambda = e^{1.707 + 1.388SST_{MDR} - 1.521SST_{TROP}} \quad (Eq.1)$$

22 where SST_{MDR} and SST_{TROP} are anomalies in the regional SST indices relative to the
23 1982–2005 average, as described in Section II.C. SST_{MDR} is the average over the

hurricane main development region (80°W - 20°W , 10°N - 25°N), and SST_{TROP} is the global, 30°S - 30°N average of SST. As discussed in Vecchi *et al.* (2011), this statistical emulator of the sensitivity of hurricane frequency to SST changes in the Zhao *et al.* (2009, 2010) high-resolution atmospheric model was trained across a broad range of climate states, including multiple realizations of the historical period and various projections of 21st century SST change. Because it was trained against a wide range of climate states, we expect that this statistical emulator should have the potential to be applicable even to future climate states. The parameters in this statistical emulator, built on the output of a high-resolution AGCM, are very similar to those that arise from modeling adjusted hurricane frequency over the 1878-2008 period (Villarini *et al.* 2012). This statistical emulator is able to reproduce much of the observed variability in hurricane activity ($r^2=0.58$; Vecchi *et al.* 2011), and its ability to recover changes in hurricane frequency compares well with hindcasts and projections from high-resolution dynamical models (*e.g.*, Zhao *et al.* 2009, 2010; Villarini *et al.* 2011a). The low computational cost of the statistical emulator allows us to efficiently perform a variety of retrospective forecasts using multiple input datasets, described below.

B. Global climate model predictions:

The statistical emulator (described above) is applied to predictions of SST from two global climate models: NOAA Geophysical Fluid Dynamics Laboratory (GFDL) CM2.1 and UKMetOffice (UKMO) Decadal Prediction System (DePreSys) Perturbed Physics Ensemble (PPE), referred to as GFDL-DecPre and UKMO-DePreSys, respectively. The forecast system specifications are summarized in Table 1. These models are just two of

1 those what will be part of the CMIP5 decadal prediction experiments, although the
2 CMIP5 version of UKMO-DePreSys is slightly different from the one used here.

3 The GFDL decadal climate hindcasts (GFDL-DecPre) are carried out over the period
4 1961-2011 using the GFDL CM2.1 coupled system (Delworth *et al.* 2006), in which both
5 the atmosphere and the ocean are initialized through a full-field assimilation to bring the
6 state of the coupled model close to observations. The initial conditions are produced with
7 the GFDL fully coupled reanalysis ECDA3.1, which is based on an ensemble Kalman
8 filter (Zhang *et al.* 2007; Zhang and Rosati 2010; Chang *et al.* 2011b) and has been
9 shown to produce a realistic ocean mean state and variability (Chang *et al.* 2012). Ten-
10 member ensembles are produced starting from the first of January every year from 1961
11 to 2011 and run for ten years. Historical radiative forcing is used for the 1961-2005
12 period and the Representative Concentration Pathways (RCP) 4.5 scenario for the
13 predictions starting after 2005. A ten-member ensemble of uninitialized runs with the
14 same forcings has also been produced to investigate the impact of initialization. This
15 forecast suite is further discussed in Rosati *et al.* (2012), and its retrospective skill in
16 predicting Atlantic Multidecadal Oscillation-like variability is described in Yang *et al.*
17 (2012).

18 DePreSys (Smith *et al.* 2007) is based on the third Hadley Centre coupled global
19 climate model, HadCM3 (Gordon *et al.* 2000). The UKMO-DePreSys Perturbed Physics
20 Ensemble (PPE; Smith *et al.* 2010) is an updated version that uses a nine-member
21 ensemble of model variants that aims to sample model uncertainties through
22 perturbations to poorly constrained atmospheric and surface parameters. Initial conditions
23 are created by relaxing the model's components toward atmospheric (European Centre for

1 Medium Range Weather Forecasting Analysis and Reanalysis) and oceanic (Smith and
2 Murphy 2007) analysis, with values assimilated as anomalies with respect to the model
3 climate. The purpose of anomaly assimilation is to minimize climate drift after the
4 assimilation is switched off, but this does not totally suppress the bias as discussed in
5 Robson (2011). The ten-year long decadal retrospective forecasts consist of nine-member
6 ensembles starting from the first of November every year from 1960 to 2005. A parallel
7 set of nine uninitialized experiments using the DePreSys-PPE is also used, and is referred
8 to as the UKMO-DePreSys uninitialized forecast runs. We use the UKMO-DePreSys-
9 PPE data, rather than the CMIP5 UKMO-DePreSys output in order to have a comparison
10 to the results of Smith *et al.* (2010); work is underway with the full suite of CMIP5
11 models (Caron *et al.* 2012, in preparation).

12 *C. Lead-dependent climatology:*

13 The statistical hurricane emulator is defined in terms of SST anomalies against the
14 1982-2005 climatology (Vecchi *et al.* 2011). The initialized and uninitialized model
15 forecasts have their own climatology, which –for initialized forecasts using both models
16 and for uninitialized forecasts using UKMO-DePreSys-PPE – can depend on the lead-
17 time of the forecast. Therefore, we define a different climatology for each experiment
18 (initialized and uninitialized), for each model (GFDL-DecPre and UKMO-DePreSys-
19 PPE). For the initialized model experiments we build a climatology that depends on lead-
20 time by averaging, for each lead-time between one and ten years, the forecasts that verify
21 in the years 1982-2005. To compute the model climatology we average all ten ensemble-
22 members for GFDL-DecPre, but since UKMO-DePreSys-PPE is a “perturbed physics
23 ensemble” a different climatology is defined for each of its nine ensemble members. Note

that a key impact of subtracting the lead-dependent climatology is to remove a systematic bias that arises in the forecasts as the models drift toward their own mean state when initialized with observations (Stockdale 1997; ICPO 2011). A key assumption is that the systematic drift of the models does not depend on initialization period – that is, that the systematic drift does not depend on the changes to the climate observing system that have occurred in the last 50 years. The assumption that the drift is stationary will be further discussed in Section IV.

D. Skill measures:

We explore two statistical measures to quantitatively assess retrospective performance: anomaly correlation coefficient (ACC), and mean squared skill score (MSSS). These statistics are not independent, but offer slightly different views of the forecast model skill. ACC is the sample correlation coefficient as a function of lead time t (or an average of lead times), between a set of forecast anomalies F'_j and observed anomalies O'_j , over $j = 1, \dots, n$ years after removing the mean of each:

$$ACC(t) = \frac{\sum_{j=1}^n (F'_j(t) \cdot O'_j(t))}{\sqrt{\sum_{j=1}^n F'^2_j(t) \sum_{j=1}^n O'^2_j(t)}} \quad (\text{Eq.2})$$

where $F'_j = F_j - \bar{F}$, $O'_j = O_j - \bar{O}$ and the overbar denotes the time mean over the climatological period 1982-2005¹, which is a function of lead time t . ACC values can range from -1 to 1, and they measure the degree to which large positive and negative excursions from the mean co-occur in the forecast and verification.

¹ The climatological period 1982-2005 is used because that is the reference period of the study that developed the statistical emulator (Vecchi *et al.* 2011).

The root-mean squared error (RMSE) is often used as a measure of accuracy of the forecasts. It is defined as the square root of the mean squared error (MSE)

$$RMSE(t) = \sqrt{MSE(t)} = \sqrt{\frac{1}{n} \sum_{j=1}^n (F'_j(t) - O'_j(t))^2} \quad (\text{Eq.3})$$

We use here a related statistical measure, the mean squared skill score (MSSS; Murphy 1998) following recommendations by Goddard *et al.* (2012). MSSS is based on the mean squared error (MSE) between the forecast and the observed climatology and represents the improvement in accuracy of the forecast over climatology:

$$MSSS(t) = 1 - \frac{MSE_F(t)}{MSE_{\bar{X}}(t)} \quad (\text{Eq.4})$$

The highest MSSS value of 1 is reached when $MSE_F = 0$ and $MSE_{\bar{X}} \neq 0$.

Instead of using climatology as reference forecast one can use the MSE of the uninitialized projections (MSE_p) to evaluate the improved skill due to initialization:

$$MSSS(t) = 1 - \frac{MSE_F(t)}{MSE_p(t)} \quad (\text{Eq.5})$$

where a positive MSSS indicates that the initialized forecasts outperform the uninitialized ones. MSSS can be expressed as a function of correlation and conditional bias (Goddard *et al.* 2012), which is useful when interpreting an improvement of skill due to initialization.

E. Assessment of statistical significance:

We explored three different estimates to assess statistical significance of the correlation results against a null of zero correlation, and to compute the confidence intervals of the retrospective correlations. For the estimates of statistical significance the effective number of degrees of freedom (N_{eff}) of the correlation of two time-series (X and

Y) was computed using the methodology described in Bretherton *et al.* (1999), using the biased estimates of autocorrelation spectrum of the various time-series:

$$N_{eff} = \frac{N}{\sum_{\tau=0}^{N-1} (1-|\tau|/N) r_{\tau}^X r_{\tau}^Y} \quad (\text{Eq.6})$$

where N is the number of samples in each time-series, and r_{τ}^X and r_{τ}^Y is the estimate of autocorrelation of each time-series at lag τ . Because of the large autocorrelation of the time-smoothed predicted and observed hurricane time-series at even long lags, the effective degrees of freedom can be considerably smaller than the number of years in the time-series. Typically, when compared to observations, the five-year mean initialized forecasts tend to have between 6-8 effective degrees of freedom and the uninitialized forecasts tend to have between 10-12 effective degrees of freedom – even though there are around fifty years of data that are compared. Without accounting for the strong autocorrelation in these time-series, one would estimate much narrower confidence intervals and a smaller p -value for the null hypothesis; failure to account for the diminished degrees of freedom can lead to a substantial overestimation of forecast skill.

Though hurricane frequency is not Normally-distributed, we are exploring multi-year averages of hurricane frequency, which allows us to approximate the distribution as Normal. To compute confidence intervals of a correlation we use a two-sided test (since it is possible that initialization could lead to degradation in performance), and use a one-sided test against the null hypothesis of zero correlation (since a significantly negative correlation would be a failure of the forecast system), we have compared the results from three methods:

- i) *Fisher's-z Transformation*: The sample estimate of the correlation coefficient between two time-series (X and Y), $r_{X,Y}$, is transformed using:

$$z_{X,Y} = 0.5 \ln \left[(1 + r_{X,Y}) / (1 - r_{X,Y}) \right] \quad (\text{Eq.7})$$

The new quantity, $z_{X,Y}$, follows a z distribution with $N_{eff}-3$ degrees of freedom (Fisher 1915, 1924; von Storch and Zwiers, 1999). Using standard z -statistic tables one can estimate the confidence intervals on the mean and test against a null of zero mean from the sample estimate, $z_{X,Y}$. To transform the confidence interval estimates of the z -statistic back to correlation space, we employ the inverse Fisher's- z Transformation:

$$r_{X,Y}^* = \frac{e^{2z_{X,Y}^*} - 1}{e^{2z_{X,Y}^*} + 1} \quad (\text{Eq.8})$$

where $z_{X,Y}^*$ is the estimate of the upper or lower bound on the confidence interval of the z -statistic and $r_{X,Y}^*$ is the estimate of the upper or lower bound on the confidence interval of the correlation coefficient.

ii) *Full distribution of the correlation coefficient:*

Johnson *et al.* (1995) provide the distribution of the sample correlation coefficient R when the population correlation coefficient ρ is equal to zero:

$$p_R(r) = \frac{\Gamma[(n-1)/2]}{\Gamma(1/2)\Gamma[(n-2)/2]} (1-r^2)^{(n-4)/2}, \quad \text{for } -1 < r < 1 \quad (\text{Eq.9})$$

where $\Gamma(\cdot)$ is the gamma function, n is the sample size. This distribution is symmetric around the zero. By using p_R , we can test the null hypothesis of no correlation at a given significance level α , by checking whether the sample correlation coefficient lies within or outside the rejection or critical region.

iii) *Monte Carlo estimate:* For sample sizes ranging between 2 and 100, we build 100,000 estimates of the distribution of the sample correlation coefficient

1 between two normally-distributed time-series of length N_{eff} and an underlying
2 correlation ρ . We sample underlying correlation coefficients between -1 and
3 1, at intervals of 0.01. From this Monte Carlo estimate of the probability
4 density function of the sample correlation coefficient, we estimate
5 significance against a null of zero correlation as the probability of a
6 correlation as large as or larger than a particular sample correlation given an
7 underlying correlation of zero. In an analogous manner, we also compute the
8 confidence intervals on the sample correlation given an underlying
9 correlation.

10 We have compared the three estimates of the confidence intervals on the
11 correlation coefficient and null test against a correlation of zero for the
12 retrospective forecast correlations, and have found that they are consistent with
13 each other. For simplicity, in the manuscript we only show the estimates from the
14 Fisher's-z transformation.

15 **III. Results**

16 *A. Retrospective Hurricane Forecasts:*

17 Figure 1 shows the five-year mean and nine-year mean (centered on the mid-point of
18 each interval) initialized and uninitialized forecasts of North Atlantic hurricane frequency
19 in GFDL-DecPre and UKMO-DePreSys-PPE compared with observations. The observed
20 record of five-year mean hurricane frequency is characterized by two distinct states with
21 low values (~ 5 -6 hurricanes per year) in the first half of the record and a shift in the mid-
22 90s (*e.g.*, Elsner et al. 2004, Li and Lund 2012) toward a more active state (~ 8 hurricanes
23 per year). The uninitialized predictions capture a tendency for an increase in hurricane

1 frequency over the late-20th century, indicating that part of the recent increase in Atlantic
2 hurricane frequency was due to changes in radiative forcing. However the uninitialized
3 experiments fail to capture the abrupt shift in the mid-1990s. The initialized retrospective
4 forecasts show better qualitative agreement to observations than do the initialized runs,
5 suggesting an improvement from initialization.

6 Despite the time averaging, both observations and the model predictions have year-to-
7 year variability in five-year North Atlantic hurricane frequency, leading to substantial
8 noise in the detection of the multi-year to decadal changes (Figure 1). The year-to-year
9 variations in the multi-year initialized forecasts are larger than that in observations, even
10 though the forecasts are ensemble averages. This result suggests that the initialized
11 forecasts have too much internal variability. An alternative interpretation, which is
12 discussed further in Section III.C below, is that the initial conditions for each model are
13 persisted too strongly, so that each year's initial climate reflects in subsequent years.

14 The anomaly correlation between the observed hurricane counts and the models
15 predictions for both initialized and uninitialized experiments is shown in Figure 2 for
16 five-year and nine-year means. A persistence forecast is given as a reference test forecast,
17 where the five-year (nine-year) mean persistence is defined as the observed average over
18 the five (nine) years that precede the model's initialization. So, for example, the
19 persistence forecast for the lead 2-6 forecast centered in 1992 (*e.g.*, initialized in 1989) is
20 the observed hurricane count averaged over 1984-1988. Consistent with Figure 1, at lead
21 2-6 the initialized retrospective predictions show higher correlations than the uninitialized
22 ones, for both models. The values are significantly different from zero and exceed the
23 values given by persistence, which is not the case for the uninitialized predictions.

1 Comparable skill is found between the two models, slightly higher in UKMO-DePreSys.
2 Computing the two-model mean increases the signal-to-noise ratio, leading to higher
3 correlations than in either individual model. At lead 2-10, all the predictions outperform
4 the persistence forecast. The decadal correlations are nominally higher in the initialized
5 retrospective predictions than in the uninitialized, with the largest values, exceeding 0.8,
6 when taking the two-model mean. This decadal skill does not come only from the first
7 few years since the correlations at lead 6 to 10 are also large (Figure 2), although the
8 improvement due to initialization is not as clear. At lead 6-10, GFDL-DecPre shows
9 larger correlations for the initialized predictions but UKMO-DePreSys indicates higher
10 values for the uninitialized runs, yielding undistinguishable values between the initialized
11 and the non-initialized experiments for the two-model mean.

12 These results suggest that coupled GCMs that account for both changes in initial state
13 and radiative forcings can lead to skillful multi-year retrospective predictions of
14 hurricane frequency. The nominal improvement due to initialization should, however, be
15 interpreted with care given the large confidence intervals associated with the point
16 estimates of the correlations (Figure 2). As discussed above in Section II.E, although the
17 observed record is 50-years long, because of the large autocorrelation of the time series
18 each year is not independent from those nearby. Hence, the effective number of degrees
19 of freedom is largely reduced to less than ten for most lead times, as indicated on Figure
20 2, based on Bretherton *et al.* (1999). Therefore, even if the initialized predictions give a
21 correlation that is statistically different from climatology and is nominally higher than in
22 the uninitialized predictions, the large confidence intervals indicate that the retrospective
23 correlation of the initialized forecasts is not different from persistence or the uninitialized

1 experiments at $p=0.1$. Some of the correlations of the initialized forecasts are
2 significantly larger than the non-initialized experiments at $p=0.2$.

3 Improvement from initialization on the two-model mean lead 2-6 forecast is close to
4 being significant even at $p=0.1$, suggesting potentially higher confidence in multi-model
5 ensembles. For the lead 2-6 and 2-10 forecasts, for both model systems there is a
6 consistent nominal improvement of retrospective correlation from initialization relative to
7 the uninitialized experiments. Because of this, and because of the small sample size, we
8 speculate that the lack of significance at $p=0.1$ may reflect a “lack of power” by the
9 significance test, rather than a “lack of effect” from initializing (Johnson 1999). For the
10 lead 6-10 forecast, however, the nominal difference between the initialized and non-
11 initialized forecasts changes sign (there is nominal indication of improvement in GFDL-
12 DecPre, but a nominal degradation in UKMO-DePreSys-PPE), so we interpret the lack of
13 significance in this case as indicating a lack of effect from initialization. Therefore, it
14 appears that the nominal improvement in the lead 2-10 forecast arises in the first part of
15 the decade, and represents potential multi-year forecast skill rather than decadal skill.

16 A lagged ensemble approach, in which past forecasts are used to augment the
17 effective ensemble size of more recent forecasts (*e.g.*, by creating a forecast where the
18 current year’s lead 1-5 and the previous year’s lead 2-6 forecasts are averaged), can lead
19 to increase in forecast performance (*e.g.*, Vecchi *et al.* (2011) showed improvement in
20 interannual hurricane forecasts from lagged ensembles). We explored the impact of
21 lagged ensembles in the retrospective hurricane forecasts (not shown) at lags of up to
22 three years (*i.e.*, averaging lead 1-5, 2-6 and 3-7 verifying the same years together)
23 resulted in nominal improvements in the correlation coefficient (on the order of 0.02-

0.05). However, the smoothing induced by the lagged ensemble led to a further reduction of degrees of freedom. Since the uncertainty in a correlation estimate increases with decreasing correlation or sample size, the uncertainty estimates on the correlation coefficient did not show substantial change: even after lagged-ensemble averaging the retrospective correlation of the uninitialized and initialized forecasts were in each other's confidence intervals.

As a complement to the skill estimate using ACC, we show in Figure 3 the MSSS for various five-year mean and nine-year mean leads. Both the improvement relative to climatology (Eq.4) and that due to initialization (Eq.5) are indicated on the x- and y-axis, respectively. None of the retrospective initialized forecasts has a negative MSSS on the x-axis, which indicates at least a nominal improvement relative to climatology. An improvement due to initialization is also suggested at all leads in GFDL-DecPre, and at most leads except 5-9 and 6-10 in UKMO-DePreSys, leading to a smaller MSSS at those lead times for the two-model mean. Both models indicate an improved skill at decadal scale due to initialization, with the highest values in UKMO-DePreSys. As shown in Goddard *et al.* (2012), the MSSS is a function of both the correlation and the conditional bias, and the higher MSSS due to initialization is mainly due to a reduction of the conditional bias that is large in the uninitialized predictions.

B) SST-source of hurricane forecast skill

Our hurricane frequency index is based on SST averaged over the tropical Atlantic and over the global tropics (Eq.1), so both quantities are potential sources for the better predictability in the initialized forecasts. We can explore retrospective forecasts and skill measures of these two indices with hope of finding the role each had in recovering the

1 past history of hurricane activity (Figure 4). Overall, there is no indication that
2 retrospective forecasts of tropical-mean SST are improved by initializing the coupled
3 GCMs (upper panels, Figure 4), with the relatively monotonic warming of the tropics
4 dominating the observed and modeled signals. The dominance of the long-term trend in
5 both SST indices cuts the effective degrees of freedom severely, to the point where for
6 tropical-mean SST interpretation of correlation as a skill metric is likely too ambiguous
7 to be useful. The GFDL-DecPre system has marginally higher retrospective correlation in
8 both SST indices than does UKMO-DePreSys, likely due to inclusion of future volcanic
9 information in its radiative forcing (Table 1). However, this nominally larger skill in
10 GFDL-DecPre for the two SST indices does not translate into even nominal increase of
11 the hurricane forecasts (Figure 2) since the volcanic signals are primarily spatially
12 uniform. Across both model systems there is a consistent nominal improvement of
13 retrospective correlation of Atlantic MDR SST predictions from initialization, but the
14 effect is small relative to the number of degrees of freedom. Only in the GFDL-DecPre
15 does the initialized forecast of MDR SST approach a significant improvement over a
16 persistence forecast. Because of the dominance of a quasi-monotonic trend, for tropical-
17 mean SST all the forecast methods (initialized and uninitialized GCM forecasts and
18 persistence) yield comparable results. For both SST indices all of the forecast
19 methodologies lead to statistically significant retrospective correlations against a null of
20 zero correlation, again largely because of the dominance of a trend.

21 The results in Figure 4 suggest that the nominal improvement in retrospective
22 correlation from initialization came from improvements to forecast of Atlantic MDR
23 SST. However, since the time series of each SST index includes a substantial component

1 that is coherent across both indices, and since the hurricane frequency emulator is based
2 on the difference between the two indices, interpreting the source of hurricane
3 predictability from each index is not necessarily straightforward, as was noted in Vecchi
4 *et al.* (2011). An alternative approach to assessing influence of each index on the role of
5 initialization on forecast skill is to use values of one index from the initialized
6 experiments and the other from the uninitialized experiments. For example, taking values
7 for SST_{MDR} from the initialized experiment, but keeping the SST_{TROP} from the
8 uninitialized one, yields comparable hurricane retrospective forecast results (Figure 5a) to
9 when both indices are taken from the initialized experiments (Figure 2). The impact of
10 initialization on SST_{MDR} yields five-year mean fluctuations of this hurricane frequency
11 index that show rather good agreement with observations for both models with a
12 correlation of 0.70 and 0.59 in GFDL-DecPre and UKMO-DePreSys, respectively (both
13 significantly different from zero correlation at $p < 0.05$) at lead 2-6. Using values for
14 SST_{MDR} from the uninitialized experiments but those of SST_{TROP} from the initialized
15 experiments leads to very different results (Fig 5.b). The correlation drops to 0.21 in
16 GFDL-DecPre and to 0.43 in UKMO-DePreSys, with neither correlation significantly
17 different from $\rho = 0$ (even at $p < 0.2$) nor either model able to reproduce the observed sharp
18 increase in the mid 90s. This indicates that the nominal improvement in correlation in the
19 initialized multi-year predictions results from a better representation of the Atlantic main
20 development region when initializing the coupled models, with little beneficial impact
21 from initialized predictions of the global mean tropical SST.

22 For the GFDL-DecPre system the difference in retrospective correlation when
23 swapping initialized/uninitialized SST_{MDR} and SST_{TROP} is significant at $p < 0.1$. Note in

Figure 5.b there is a large increase in hurricane frequency around 2005 in GFDL-DecPre, as appeared in Fig1.a. This increase, which we currently consider to be spurious, is a large contributor to the reduction in correlation from the impact of initialization on tropical-mean SST in the GFDL model. There is a coincidence between the global implementation of the “Array for Real time Geostrophic Oceanography” (or Argo) drifting float profiles in 2003 and the spurious shift of nine-year forecasts centered around 2005-2006, suggesting that enhanced observational sampling after 2003 may have led to a change in the lead-dependent climatology. Experiments are underway to test this possibility. The lack of such a spurious increase in UKMO-DePreSys could arise from different initialization processes, or from the fact that the last initialized forecast in UKMO-DePreSys begins in 2006 – so the late spike would not be evident. Were the introduction of Argo found to be the driver of this spurious increase, in addition to developing methods to minimize the impact of observing system changes, the impact of other large changes to the observing system must also be explored (*e.g.*, the introduction of altimetry in the early 1990s and the completion of the TAO array in the mid-1990s).

C) Role of the mid-1990s climate shift:

The nominal improvement in skill due to initialization should be interpreted with care. Even if the initialized retrospective predictions outperform climatology at almost all lead times (Figure 3), the skill could still come from persistence – just persistence that cannot be captured with our observationally-based persistence model. Figure 6a and 6b compare the retrospective predictions of hurricane frequency for five-year means ranging between lead 1-6 to lead 6-10. The forecasts at each lead show a tendency to have a systematic one year shift with respect to the preceding lead, with the mid-1990s shift in

1 each model trailing in time for longer leads rather than capturing the observed 1995 shift
2 (*e.g.*, Elsner et al. 2004, Li and Lund 2012) at the right time. By performing change point
3 analysis, Pettitt test, to the models' retrospective predictions, we find a shift in forecasts
4 initialized in 1991 in UKMO-DePreSys and forecasts initialized in 1995 in GFDL-
5 DecPre. This tendency for forecasts to lock across the shift can be seen more clearly
6 when the same time series are plotted as a function of initialization year instead of
7 verification time (Fig 6c and 6d): forecasts initialized the same year are very similar to
8 each other, independent of when they verify. Notice that the mid-90s shift for each model
9 appears at the same initialization year for all lead times, as does the potentially spurious
10 mid-2000s shift in GFDL-DecPre.

11 Up to now we have been largely comparing the results of forecasts initialized
12 different years at the same lead, without focusing on the evolution of hurricane counts of
13 each forecast as the lead increases. A correct forecast of the mid-1990s climate shift
14 would have indicated at some point prior to the shift that there was an increased
15 probability of hurricane frequency increasing in time. For example, if a forecast
16 initialized in early 1991 showed counts averaged in 1992-1996 that were larger than
17 those in 1991, or an increased number of ensemble members with large increases, one
18 would have evidence for a future shift. Do these two forecast systems produce such a
19 shift? Figure 7 shows that in the observational record, reflecting the rapid increase in
20 frequency in 1995, the difference in hurricane counts averaged over the five years
21 following the years 1991 through 1994 exceeded the counts over each of those years by
22 an unusually large amount, relative to the distribution over the 1961-2006 period.
23 However, neither forecast system (colored lines in Figure 7) shows a tendency for their

1 forecasts to increase in time relative to the first forecast year when initialized in the early
2 1990s. In fact, there is a nominal tendency for these forecasts to decrease in time from the
3 first forecast year, relative to the distribution of tendencies across all initialization dates,
4 1961-2006. That is, the models did not forecast a *tendency* towards higher frequency in
5 the mid-1990s (Figure 7), even though the sequence of forecast *values* exhibits a climate
6 shift in the mid-1990s (Figures 1, 6).

7 To further highlight the influence of the mid-1990s shift on the retrospective skill
8 estimation, we explore forecast performance after removing the mid-90s shift from both
9 the forecasts and the observations. The shift is “removed” by simply referencing each
10 period before and after the 1994-1995 shift to its own climatology; for instance, the time-
11 mean hurricane count preceding 1995 is removed from all years before 1995, and the
12 time-mean hurricane count following 1995 is removed from all years after 1995. We note
13 that using each model’s change-point instead of 1995 does not affect the character of the
14 results. Figures 8 and 9 indicate that removing the shift leads to a substantial reduction of
15 correlation in the initialized predictions at lead 2-6 (particularly for UKMO-DePreSys-
16 PPE), and no indication of skill beyond that lead time, further confirming that the decadal
17 signal is dominated by the trend that arises from the existence of the mid-90s change
18 point. Therefore, future real (as opposed to the retrospective forecasts explored here)
19 multi-year and decadal predictions of hurricane frequency should not be expected to
20 show the same skill as over the 1961-2011 period unless there are change points of
21 similar character to the mid-1990s shift. Our results are encouraging for the feasibility of
22 multi-year forecasts of hurricane frequency with the current prediction systems.
23 However, this analysis highlights that substantial challenges remain – or, viewed more

1 optimistically, that it is possible to improve the performance of the system beyond its
2 current capability.

3 An interesting side effect of removing the mid-1990s shift is to increase the effective
4 degrees of freedom, narrowing the confidence intervals associated with the point
5 estimates of the correlation coefficient (compare Figures 2 and 9). In addition, the
6 retrospective correlation in the uninitialized forecasts without change-point disappeared –
7 since it largely arose from the projection of the observed shift onto the models' forced
8 trend over this period. In this modified context, there is now indication that for the GFDL
9 model and the two-model ensemble the correlations (although lower than in the case
10 including the shift; Figure 2) are significantly higher than those of the uninitialized
11 versions of the model at lead 2-6. That is, there is significant (at $p<0.1$) indication that
12 GFDL-DecPre and the two-model ensemble may be able to predict the types of variations
13 in hurricane frequency that occurred in the early-1980s and early-1990s better than the
14 uninitialized experiments. In Figure 2, the nominal improvement from initialization in the
15 correlation of the lead 2-6 and lead 6-10 mean hurricane counts in GFDL-CM2.1 was
16 larger than that for the lead 2-10 forecasts; this may reflect the ability of GFDL-CM2.1 to
17 retrospectively forecast some multi-year variations beyond the 1994-1995 climate shift –
18 which is the dominant signal in the nine-year running counts. This further highlights the
19 limitations of a data record that is short relative to the dominant timescales in order to
20 assess the impact of multi-year forecast skill. While it is entirely possible that some of the
21 non-significant differences between the initialized and uninitialized models shown in
22 Figures 2 and 3 could become significant from a longer record, it is also possible that the
23 impact of initialization could also decrease and remain non-significant in a longer record.

1 **IV Summary and Discussion**

2 The predictability of North Atlantic hurricane frequency has been investigated in
3 two global coupled models initialized towards estimates of the observed climate state.
4 We find encouraging performance of initialized retrospective forecasts of hurricane
5 activity on multi-year to decadal timescales, with statistically significant (against a null of
6 zero) retrospective correlation at leads 2-6, 6-10 and 2-10 years when accounting for both
7 initialization and radiative forcing changes. However, although there is a nominal
8 increase in correlation in the initialized forecasts, relative to the uninitialized forecasts, is
9 not statistically significant because of the small independent sample size even with over
10 50 years of data. The two systems explored, GFDL-DecPre and UKMO-DePreSys-PPE,
11 show comparable skill, with the best results obtained when using the two-model mean.
12 These results are encouraging but need to be interpreted with care because of the short
13 observational record and the persistent character of the time series we are trying to
14 predict: the confidence intervals associated with all the forecasts are large and the
15 difference between the initialized and the uninitialized forecasts is not statistically
16 significant at $p=0.1$. Using two-model and lagged averages leads to nominal increases in
17 the correlation, which encourages the pursuit of broader multi-model studies (e.g., Caron
18 *et al.* 2012, in preparation).

19 The observed time series of North Atlantic hurricane frequency is dominated by a
20 strong and abrupt rise in 1995 leading to a trend over the 1961-2011 period. The high
21 correlations of the retrospective predictions of North Atlantic hurricane frequency depend
22 on the presence of this shift. While both models show an increase in hurricane frequency
23 in the mid-90s, leading to high correlations with observations, the increase is not actually

1 predicted by the evolution of the models, but is present in the initial state (*i.e.*, forecasts
2 initialized after the shift exhibited by each model remain high, but those initialized prior
3 do not show the shift). This reduces our confidence that a similar shift in a near future
4 could be successfully predicted with the current prediction systems. It also highlights the
5 need to better understand the origin of the change point in the observations and assess
6 whether the modeled mechanisms are consistent with those in the real world (*e.g.*,
7 Robson *et al.* 2012). Relatedly, although our results suggest that the initialized forecasts
8 outperform a persistence forecast built from the observational record, we do not exclude
9 persistence as a source of skill. The initialized retrospective predictions persist the state in
10 which they are initialized, leading each model to produce a shift in the mid-90s. This long
11 persistence in the hurricane index was not expected because, unlike the subpolar North
12 Atlantic where the ocean imparts a strong inertia to the system, the hurricane index used
13 in this study is built using SSTs in the tropical Atlantic and globally, which exhibit
14 substantial variations from year to year. The mechanisms behind this persistence in the
15 hurricane index are under exploration.

16 The highest values of hurricane frequency in the observed record appear around 2005,
17 but are not predicted by either model. GFDL-DecPre shows a comparable rise but five to
18 ten years later than observed, whereas UKMO-DePreSys shows a more modest increase
19 with a several-year delay as well. This period coincides with a fundamental change in the
20 ocean observing system, with the global introduction of Argo floats after 2003 bringing a
21 considerably better coverage of the surface and subsurface ocean. This data discontinuity
22 could present problems in the forecasts, as we expect the lead-dependent climatology to
23 depend on the difference between the state to which a model is initialized and that to

1 which it naturally tends. Forecasts with GFDL-DecPre that extend past the present
2 suggest an increase in hurricane frequency through the mid-2010s (Figure 1). However,
3 observations have been tending in the opposite direction, with recent years being less
4 active than those in the mid-2000s. We are skeptical of these predictions towards an even
5 more active state in the next ten years, in part because of the unusual role of tropical-
6 mean SSTs in driving this predicted increase (Figure 5). Our current hypothesis is that
7 this increase predicted by the statistical-dynamical hybrid system with GFDL-DecPre is
8 spurious, and reflects the impact of increased data density in the 2000s on the GFDL-
9 DecPre drift. Experiments are underway to test this hypothesis. If this hypothesis is
10 correct, then a more plausible prediction for the coming years is that shown in the left
11 panel of Figure 5, in which there is a tendency towards a reduction of hurricane
12 frequency in coming years. Interpretation of these forecasts needs to be keenly
13 constrained by our knowledge of changing observing practices both in the predictands
14 (*e.g.*, Vecchi and Knutson 2008, 2011; Landsea *et al.* 2010; Villarini *et al.* 2011b) and in
15 the observations used to initialize the climate model (*e.g.*, Zhang *et al.* 2007).

16 Figures 6 and 7 show that these models did not dynamically predict the mid-1990s
17 shift in hurricane frequency, but rather persist shifts that exist in the initial conditions for
18 several years. That is, the large correlations and MSSS values shown in Figures 2 and 3
19 do not come from predicting the dynamical evolution of the climate system that leads to
20 the shifts in hurricane frequency, but from “recognizing” that a climate shift has occurred
21 and persisting that shift.

22 Identifying the source of skill in retrospective predictions is key to the success of
23 future forecasts. Recent studies (Mann and Emanuel 2007; Evan *et al.* 2009; Smith *et al.*

2010; Villarini and Vecchi 2012b,c) have argued that the recent (since the 1980s) increase of Atlantic hurricane activity was not caused by internal variability alone but also included an externally-forced component driven largely by changing aerosol concentrations. Our results partially support this interpretation, indicating high correlations that are significantly different from zero at lead 2-10 in the uninitialized forecasts, suggesting the external forcings as partly driving the decadal fluctuations of Atlantic hurricane frequency in the two models. However, the sharp mid-90s rise that dominates the decadal variations of Atlantic hurricane frequency in observations is not retrospectively predicted in the uninitialized experiments. Its better representation in the initialized predictions could be interpreted as an indication of a key role for internal variability in the mid-1990s shift, supporting various studies (*e.g.*, Zhang and Delworth 2005,2006,2009; Robson *et al.* 2012; Yeager *et al.* 2012; Msadek *et al.* in preparation). However, since the models did not successfully predict the trajectory of the mid-90s hurricane shift, the nominal improvement from initialization could also reflect a failure in the radiative forcing/response in these models that is corrected when they are constrained with observations. A key use of these multi-year forecast experiments will be to help constrain our interpretation of past climate changes.

Our results indicate that the impact of initialization on forecasts of the Atlantic main development region (MDR) was key to the higher skill in the initialized forecasts (Figures 4 and 5). However, this does not exclude a remote influence on Atlantic tropical storms. As discussed above, the impact of initialization on the remote tropics in the GFDL model led to a dramatic increase in the predictions after 2003; we speculate that this dramatic increase is an artifact of changing observing practices, and experiments are

underway to test this hypothesis and build mechanisms to correct this effect. Zhang and Delworth (2006) suggested that multi-year changes in hurricane activity could be driven by changes to the heat-transport over the entire North Atlantic. Smith *et al.* (2010) and Dunstone *et al.* (2011) further suggested that the subpolar North Atlantic was the main source of multi-year predictability of Atlantic hurricane frequency and showed that initializing that region lead to predictability of tropical storms in their model. The North Atlantic also stands out as the region where initialized forecasts outperform uninitialized ones in the GFDL model (Rosati *et al.* 2012; Yang *et al.* 2012; Msadek *et al.* in preparation), suggesting a link between North Atlantic variability and Atlantic tropical storm predictability in GFDL CM2.1. Meanwhile, Kang *et al.* (2008) showed that changes in the North Atlantic could lead to changes in atmospheric circulation over the tropical Atlantic in GFDL CM2.1. However, in our retrospective forecasts of hurricane activity, the relevant source of skill must have been present in tropical Atlantic SST – so any role for extratropical forcing must have involved a subsequent change to tropical Atlantic SST.

Our results show retrospective correlations that are comparable with those described in Smith *et al.* (2010), indicating initializing a climate model and accounting for radiative forcing changes, together, can lead to significant retrospective skill in multi-year initialized (relative to a null hypothesis of zero correlation). However, in contrast to Smith *et al.* (2010), we do not find statistically significant evidence for an enhancement of skill from initialization: there is encouraging nominal improvement, but the record is too short and the time series too autocorrelated to reject the null of no change from initialization. Because of the consistency between the two model systems and the visual

1 improvement, we hypothesize that for the lead 2-6 and lead 2-10 forecasts the lack of
2 significance of the improvement from initialization may be an indication of lack of power
3 by the statistical test (arising from too few degrees of freedom) rather than a lack of effect
4 of initialization. Therefore, additional years could lead to enhancement of our confidence,
5 since there are currently only ~ 7 effective degrees of freedom when comparing these
6 forecasts to observations. However, the large autocorrelation of the time series indicates
7 that we require about seven years of data to gain a degree of freedom – so many years
8 will be required to improve our confidence, even if we include the past 50 years in future
9 estimates of forecast skill.

10 Despite high correlation values, the mean retrospective skill of these forecasts
11 may provide a poor and even misleading guide to the future performance of a decadal
12 prediction system if not interpreted with care. In the absence of a major climate shift, like
13 the 1994-1995 shift, or in the absence of ability to predict these shift, the long-term
14 estimates of correlation (*e.g.*, 0.6-0.9) are not representative, and the lower retrospective
15 correlations assessed after removing the shift (*e.g.*, 0-0.4, as shown in Figures 6 and 7)
16 may be closer to those one should expect. Further, it is important to recall that the
17 initialized forecasts succeed in describing the existence of the 1994-1995 shift in the
18 number of hurricanes across forecasts at a given lead, but they failed to capture the
19 trajectory towards an increase within forecasts initialized in the early-1990s. Further, the
20 potential impact of changing observing systems through the introduction of the Argo
21 array may continue in the coming years. In particular, this potential limitation needs to be
22 used in interpreting the GFDL-CM2.1 system's forecast for increased hurricane
23 frequency over the next few years, since it appears to be driven primarily by a shift in the

forecasts of non-Atlantic tropical SSTs after the introduction of Argo, which is consistent with the impact of a modified model “drift” when previously poorly-sampled areas of the globe are constrained more closely to observations. Experiments are underway to test the impact of changing observing systems of multi-year forecast performance, and to build techniques to overcome inhomogeneities such as these.

Acknowledgments:

We are grateful to Doug Smith (UK Met Office) for making the UKMO-DePreSys PPE data available. We thank Ming Zhao and Tom Knutson for comments and suggestions.

References:

- Alessandri, A., A. Borrelli, S. Gualdi, E. Scoccimarro, and S. Masina, Tropical cyclone count forecasting using a dynamical seasonal prediction system: Sensitivity to improved ocean initialization, *Journal of Climate*, **24**, 2963-2982, 2011.
- Bender, M.A., T.R. Knutson, R.E. Tuleya, J.J. Sirutis, G.A. Vecchi, S.T. Garner, and I.M. Held, 2010: Model impact of anthropogenic warming on the frequency of intense Atlantic hurricanes. *Science* **327**, 454–458.
- Booth, B.B., N.J. Dunstone, P.R. Halloran, T. Andrews, and N. Bellouin, 2012: Aerosols implicated as a prime driver of twentieth-century North Atlantic climate variability. *Nature*, **484**, 228-232.
- Bretherton, C.S., M. Widmann, V.P. Dymnikov, J.M. Wallace, and I. Bladé, 1999: The Effective number of spatial degrees of freedom of a time-varying field. *J. Climate*, **12**, 1990-2009.

- 1 Broccoli, A.J., and S. Manabe, 1990: Can existing climate models be used to study
2 anthropogenic changes in tropical cyclone climate? *Geophys. Res. Lett.*, **17**, 1917-
3 1920.
- 4 Camargo, S.J., A.G. Barnston, P. Klotzbach, and C.W. Landsea, 2007a: *Seasonal tropical*
5 *cyclone forecasts*, World Meteorological Organization Bulletin, 56, 297-309.
- 6 ———, K.A. Emanuel, and A.H. Sobel, 2007b: Use of a genesis potential index to
7 diagnose ENSO effects on tropical cyclone genesis. *J. Climate*, **20**, 4819–4834.
- 8 ———, M. Ting, and Y. Kushnir, 2012: Influence of local and remote SST on North
9 Atlantic tropical cyclone potential intensity. *Climate Dynamics (submitted)*.
- 10 Chang, C.-Y., J.C.H. Chiang, M.F. Wehner, A. Friedman, and R. Ruedy, 2011a: Sulfate
11 aerosol control of tropical Atlantic climate over the 20th century. *Journal of Climate*,
12 **24**, 2540–2555.
- 13 Chang, Y.-S., S. Zhang, and A. Rosati, 2011b: Improvement of salinity representation in
14 an ensemble coupled data assimilation system using pseudo salinity profiles.
15 *Geophysical Research Letters*, **38**, L13609, DOI:10.1029/2011GL048064.
- 16 Chang, Y.-S., S. Zhang, A. Rosati, T. Delworth, and W. F. Stern, 2012: An assessment of
17 oceanic variability for 1960-2010 from the GFDL ensemble coupled data
18 assimilation, *Climate Dynamic* (in press).
- 19 Chen, J.H., and S.J. Lin, 2011: The remarkable predictability of inter-annual variability
20 of Atlantic hurricanes during the past decade. *Geophysical Research Letters*, **38**
21 (L11804), doi:10.1029/2011GL047629.
- 22 Collins, M., et al., 2006: Interannual to decadal climate predictability in the North
23 Atlantic: A multimodel-ensemble study, *J. Climate*, **19**, 1195–1203.

- 1 Delworth, T. L., and Coauthors, 2006: GFDL's CM2 global coupled climate models. Part
2 I: Formulation and simulation characteristics, *Journal of Climate*, **19**, 643-674.
- 3 ———, and Dixon K.W., 2006: Have anthropogenic aerosols delayed a greenhouse gas-
4 induced weakening of the North Atlantic thermohaline circulation? *Geophysical*
5 *Research Letters*, **33**, L02606, DOI:10.1029/2005GL024980.
- 6 Elsner, J.B., and T.H. Jagger, 2006: Prediction models for annual U.S. hurricane counts,
7 *Journal of Climate*, **19**, 2935-2952.
- 8 ———, X. Niu, and T.H. Jagger, 2004: Detecting shifts in hurricane rates using a Markov
9 Chain Monte Carlo approach, *Journal of Climate*, **17**, 2652–2666.
- 10 Emanuel, K. A., 1987: The dependence of hurricane intensity on climate. *Nature* **326**,
11 483–485.
- 12 ———, Increasing destructiveness of tropical cyclones over the past 30 years. *Nature*, **436**,
13 686–688, 2005.
- 14 ———, 2007: Environmental factors affecting tropical cyclone power dissipation. *J. Clim.*
15 **20**, 5497–5509.
- 16 ———, R. Sundararajan, and J. Williams, Hurricanes and global warming—Results from
17 downscaling IPCC AR4 simulations. *Bull. Amer. Meteor. Soc.*, **89**, 347–367,
18 2008.
- 19 Evan, A.T., D.J. Vimont, A.K. Heidinger, J.P. Kossin, and R. Bennartz, 2009: The role of
20 aerosols in the evolution of tropical North Atlantic Ocean temperature anomalies.
21 *Science*, **324**, 778–781.
- 22 Fisher, R.A., 1915: Frequency distribution of the values of the correlation coefficient in
23 samples from an indefinitely large population. *Biometrika*, **10**, 507-521.

1 Fisher, R.A., 1924: The distribution of the partial correlation coefficient. *Metron*, **3**, 329-
2 332.

3 Goddard L., A. Kumar, A. Solomon, D. Smith, G. Boer, P. Gonzalez, V. Kharin, W.
4 Merryfield, C. Deser, S. Mason, B. Kirtman, R. Msadek, R. Sutton, E. Hawkins,
5 T. Fricker, G. Hegerl, C. Ferro, D. Stephenson, G.A. Meehl, T. Stockdale, R.
6 Burgman, A. Greene, Y. Kushnir, M. Newman, J. Carton, I. Fukumori, T.
7 Delworth, 2012: A verification framework for interannual-to-decadal predictions
8 experiments, *Climate Dynamics*, under revision

9 Gordon, C., C. Cooper, C. Senior, H. Banks, J. Gregory, T. Johns, J. Mitchell, and R.
10 Wood, 2000: The simulation of SST, sea ice extents and ocean heat transports in a
11 version of the Hadley Centre coupled model without flux adjustments, *Climate*
12 *Dynamics*, **16**, 147–168.

13 Gray, W.M., 1984: Atlantic seasonal hurricane frequency. Part I: El Niño and 30 mb
14 quasi-biennial oscillation influences, *Monthly Weather Review*, **112**, 1649-1668.

15 Griffies, S.M., and K. Bryan, 1997a: Predictability of North Atlantic multidecadal
16 climate variability, *Science*, **275**(5297), 181.

17 ———, and K. Bryan, 1997b: A predictability study of simulated North Atlantic
18 multidecadal variability, *Climate Dynamics*, **13**, 459–487

19 Gualdi, S., E. Scoccimarro, and A. Navarra, 2008: Changes in tropical cyclone activity
20 due to global warming: Results from a high-resolution coupled general circulation
21 model. *J. Climate*, **21**, 5204–5228.

22 Hawkins, E., and R. Sutton, 2009. The potential to narrow uncertainty in regional climate
23 predictions. *Bulletin of the American Meteorological Society*, **90**, 1095–1107.

1 ICPO (International CLIVAR Project Office), 2011: Decadal and bias correction for
2 decadal climate predictions. January. International CLIVAR Project Office,
3 CLIVAR Publication Series No.150, 6pp. Available from
4 http://eprints.soton.ac.uk/171975/1/150_Bias_Correction.pdf

5 Jarvinen, B.R., C.J. Neumann, and M.A.S. Davis, 1984: A tropical cyclone data tape for
6 the North Atlantic Basin, 1886–1983: Contents, limitations, and uses. Tech.
7 Memo. NWS NHC 22, National Oceanic and Atmospheric Administration, 24 pp.

8 Johnson, D.H., 1999: The insignificance of significance testing. *J. of Wildlife*
9 *Management*, **63**(3), 763-772.

10 Johnson, N.L., S. Kotz, and N. Balakrishnan, 1995: *Continuous Univariate Distributions*
11 (volume 2), Wiley, 752 pages.

12 Kalnay and coauthors, The NCEP/NCAR 40-year reanalysis project. *Bull. Amer.*
13 *Meteorol. Soc.*, 77(3), 437-471, 1996.

14 Kim, H.-M., and P.J. Webster, Extended-range seasonal hurricane forecasts for the North
15 Atlantic with a hybrid dynamical-statistical model, *Geophys. Res. Lett.*, **37**,
16 L21705, doi:10.1029/2010GL044792, 2010.

17 Klotzbach, P.J., and W.M. Gray, 2009: Twenty-five years of Atlantic basin seasonal
18 hurricane forecasts, *Geophysical Research Letters*, **36** (L09711),
19 doi:10.1029/2009GL037580.

20 Knight, J.R., R.J. Allan, C.K. Folland, M. Vellinga, and M.E. Mann, 2005: A signature of
21 persistent natural thermohaline circulation cycles in observed climate. *Geophys.*
22 *Res. Lett.*, **32**, L20708, doi:10.1029/2005GL024233.

- 1 Knutson, T.R., J.J. Sirutis, S.T. Garner, I. Held, and R.E. Tuleya, 2007: Simulation of
- 2 recent increase of Atlantic hurricane activity using an 18-km-grid regional model.
- 3 *Bull. Amer. Meteor. Soc.*, **88**, 1549–1565.
- 4 ———, ———, ———, G.A. Vecchi, and I. Held, 2008: Simulated reduction in Atlantic
- 5 hurricane frequency under twenty-first- century warming conditions. *Nat. Geosci.*,
- 6 **1**(6), 359–364.
- 7 ———, *et al.*, Tropical cyclones and climate change. *Nature Geoscience* **3**, 157–163, 2010.
- 8 Landsea, C.W., G.A. Vecchi, L. Bengtsson, and T.R. Knutson, 2009: Impact of Duration
- 9 Thresholds on Atlantic Tropical Cyclone Counts. *J. Climate*, **23**, 2508-2519
- 10 LaRow, T. E., Y. K. Lim, D. W. Shin, E. P. Chassignet, and S. Cocke, 2008: Atlantic
- 11 basin seasonal hurricane simulations. *J. Climate*, **21**, 3191–3206.
- 12 ———, L. Stefanova, D. W. Shin, and S. Cocke, Seasonal Atlantic tropical cyclone
- 13 hindcasting/forecasting using two sea surface temperature datasets, *Geophysical*
- 14 *Research Letters*, **37**, 1-5, doi:10.1029/2009GL041459, 2010.
- 15 Latif, M., N. Keenlyside, and J. Bader, 2007: Tropical sea surface temperature, vertical
- 16 wind shear, and hurricane development. *Geophysical Research Letters*, **34**,
- 17 L01710, doi:10.1029/2006GL027969.
- 18 Li, S., and R. Lund, Multiple changepoint detection via genetic algorithms, 2012: *J.*
- 19 *Climate*, **25**, 674-686.
- 20 MacAdie, C.J., C.W. Landsea, C.J. Neumann, J.E. David, E. Blake, and G.R. Hammer,
- 21 2009: *Tropical cyclones of the North Atlantic Ocean, 1851-2006*, Technical
- 22 Memo, National Climatic Data Center in cooperation with the TCP/National
- 23 Hurricane Center.

1 Mann, M.E., and K.A. Emanuel, 2006: Atlantic hurricane trends linked to climate
2 change. *Eos, Transactions of the American Geophysical Union*, **87**,
3 doi:10.1029/2006EO240001.

4 Mendelsohn, R., K. Emanuel, S. Chonabayashi, and L. Bakkensen, 2012: The impact of
5 climate change on global tropical cyclone damage, *Nature Climate Change*, **2**, 205-
6 209.

7 van Oldenborgh, G. J. and Doblas-Reyes, F. J. and Wouters, B. and Hazeleger, W., 2011:
8 Decadal prediction skill in a multi-model ensemble. *Climate Dynamics*, **38**, 1263-
9 1280.

10 Oouchi, K., J. Yoshimura, H. Yoshimura, R. Mizuta, S. Kusumoki, and A. Noda, 2006:
11 Tropical cyclone climatology in a global warming climate as simulated in a 20-km-
12 mesh global atmospheric model: Frequency and wind intensity analysis. *Journal of*
13 *the Meteorological Society of Japan* **84**, 259–276.

14 Peduzzi, P., B. Chatenoux, H. Dao, A. De Bono, C. Herold, J. Kossin, F. Mouton, and O.
15 Nordbeck, Global trends in tropical cyclone risk, *Nature Climate Change*, **2**, 289-
16 294, 2012.

17 Pielke, R. A. Jr and coauthors, Normalized hurricane damages in the United States:
18 1900–2005 *Nat. Hazard. Rev.*, **9**, 29–42, 2008.

19 Pohlmann, H., M. Botzet, M. Latif, A. Roesch, M. Wild, and P. Tschuck, 2004:
20 Estimating the decadal predictability of a coupled AOGCM, *J. Climate*, **17**(22),
21 4463–4472.

1 Ramsay, H. A., and A. H. Sobel, 2011: Effects of relative and absolute sea surface
2 temperature on tropical cyclone potential intensity using a single-column model.
3 *J. Climate*, **24**, 183–193.

4 Rayner, N.A., D.E. Parker, E.B. Horton, C.K. Folland, L.V. Alexander, D.P. Rowell,
5 E.C. Kent, and A. Kaplan, 2003: Global analyses of sea surface temperature, sea
6 ice, and night marine air temperature since the late nineteenth century. *J.*
7 *Geophys. Res.*, **108**, 4407, doi:10.1029/2002JD002670.

8 Robson, J., 2011: Understanding the performance of a decadal prediction system. U.
9 Reading Ph.D. Thesis, available at:
10 http://www.met.reading.ac.uk/~swr06jir/thesis/JIR_thesis.pdf

11 ———, R. Sutton, K. Lohmann, D. Smith, and M. Palmer, 2012: Causes of the rapid
12 warming of the North Atlantic Ocean in the mid 1990s. *Journal of Climate*, **25**,
13 4116–4134.

14 Rosati, A. and co-authors, 2012: Decadal Climate Prediction Experiments at GFDL. *J.*
15 *Climate* (submitted).

16 Rotstayn, L. D., U Lohmann, 2002: Tropical Rainfall Trends and the Indirect Aerosol
17 Effect. *J. Climate*, **15**, 2103–2116. doi: 10.1175/1520-0442.

18 Shen, W., R. E. Tuleya, and I. Ginis, 2000: A sensitivity study of the thermodynamic
19 environment on GFDL model hurricane intensity: Implications for global
20 warming, *Journal of Climate*, **13**, 109–121.

21 Smith, D. M., Smith, D., and J. Murphy, 2007: An objective ocean temperature and
22 salinity analysis using covariances from a global climate model, *Journal of*
23 *Geophysical Research*, **112**, doi:10.1029/2005JC003172.

- 1 ———, S. Cusack, A. W. Colman, C. K. Folland, G. R. Harris, and J. M. Murphy 2007:
2 Improved Surface Temperature Prediction for the Coming Decade from a Global
3 Climate Model, *Science*, **317**, 796–799.
- 4 ———, R. Eade, N.J. Dunstone, D. Fereday, J. M. Murphy, H. Pohlmann, and A.A. Scaife,
5 2010: Skillful multi-year predictions of Atlantic hurricane frequency, *Nature*
6 *Geoscience*, **3**, 846-849.
- 7 Smith, T.M., R.W. Reynolds, T.C. Peterson, and J. Lawrimore, 2008: Improvement to
8 NOAA’s historical merged land–ocean surface temperature analysis (1880–2006).
9 *J. Climate*, **21**, 2283–2296.
- 10 Sobel, A.H., I.M. Held, and C.S. Bretherton, 2002: The ENSO signal in tropical
11 tropospheric temperature. *J. Climate*, **15**, 2702–2706.
- 12 Stockdale, T.N., 1997: Coupled ocean–atmosphere forecasts in the presence of climate
13 drift, *Mon. Wea. Rev.*, **125**, 809–818.
- 14 Sugi, M., H. Murakami, and J. Yoshimura, 2009: A reduction in global tropical cyclone
15 frequency due to global warming. *SOLA*, **5**, 164–167.
- 16 ———, ———, and ———, 2012: On the mechanism of tropical cyclone frequency changes
17 due to global warming. *J. Meteorol. Soc. Japan*, **90A**, 397-408.
- 18 Sutton, R.T. and D.L.R. Hodson, 2005: Atlantic Ocean forcing of North American and
19 European summer climate, *Science*, **309**(5731), 115-118.
- 20 Swanson, K.L., 2008: Nonlocality of Atlantic tropical cyclone intensities. *Geochemistry*
21 *Geophysics Geosystems* **9**, Q04V01, doi:10.1029/ 2007GC00184.

1 Tang, B.H., and J.D. Neelin, 2004: ENSO influence on Atlantic hurricanes via
2 tropospheric warming, *Geophysical Research Letters*, **31** (L24204),
3 doi:10.1029/2004GL021072.

4 Taylor, K.E., R.J. Stouffer, and G.A. Meehl, 2012: An overview of CMIP5 and the
5 experiment design. *Bulletin of the American Meteorological Society*, **93**, 485-498.

6 Vecchi, G.A., and B.J. Soden, 2007a: Effect of remote sea surface temperature change on
7 tropical cyclone potential intensity. *Nature*, **450**, 1066–1071.

8 ———, and ———, 2007b: Global warming and the weakening of the tropical circulation. *J.*
9 *Climate*, **20**(17), 4316-4340.

10 ———, and T.R. Knutson, 2008: On estimates of historical North Atlantic tropical cyclone
11 activity. *J. Climate*, **21**(14), 3580-3600.

12 ———, and ———, 2011: Estimating annual numbers of Atlantic hurricanes missing from
13 the HURDAT database (1878-1965) using ship track density. *J. Climate*, **24**(6),
14 1736-1746

15 ———, K.L. Swanson, and B.J. Soden, 2008: Whither Hurricane Activity? *Science* **322**
16 (5902), 687-689.

17 ———, M. Zhao, H. Wang, G. Villarini, A. Rosati, A. Kumar, I. M. Held, and R. Gudgel,
18 2011: Statistical-dynamical predictions of seasonal North Atlantic hurricane
19 activity, *Monthly Weather Review*, **139**(4), 1070-1082.

20 Villarini, G., and G.A. Vecchi. 2012a: North Atlantic Power Dissipation Index (PDI) and
21 Accumulated Cyclone Energy (ACE): Statistical modeling and sensitivity to sea
22 surface temperature changes. *Journal of Climate* **25**(2), 625-637.

- , and ——, 2012b: Twenty-first-century projections of North Atlantic tropical storms from CMIP5 models, *Nature Climate Change*, doi:10.1038/NCLIMATE1530.
- , and ——, 2012c: Projected increases in North Atlantic tropical cyclone intensity from CMIP5 models, submitted to *Journal of Climate*.
- , and ——, 2012d: Multi-season lead forecast of the North Atlantic Power Dissipation Index (PDI) and Accumulated Cyclone Energy (ACE). Submitted to *J. Climate*.
- , ——, and J.A. Smith, 2010: Modeling of the dependence of tropical storm counts in the North Atlantic Basin on climate indices. *Monthly Weather Review* **138**(7), 2681–2705.
- , ——, and ——, 2012: U.S. landfalling and North Atlantic hurricanes: Statistical modeling of their frequencies and ratios. *Monthly Weather Review*, 140, 44–65.
- , ——, T.R. Knutson, M. Zhao and J.A. Smith, 2011a: Reconciling differing model projections of changes in the frequency of tropical storms in the North Atlantic basin in a warmer climate, *J. Climate*, **24**(13), 3224–3238.
- , ——, ——, and J.A. Smith, 2011b: Is the recorded increase in short duration North Atlantic tropical storms spurious? *J. Geophys. Res.* **116**, D10114, doi:10.1029/2010JD015493.
- Vitart, F., Seasonal forecasting of tropical storm frequency using a multi-model ensemble, *Quarterly Journal of the Royal Meteorological Society*, **132**, 647–666, 2006.

- 1 —, M. Huddleston, D. Deque, T. Palmer, T. Stockdale, M. Davey, S. Ineson, and
2 A. Weisheimer, 2007: Dynamically-based seasonal forecasts of Atlantic tropical
3 storm activity issued in June by EUROSIP, *Geophysical Research Letters*, **34**
4 (L16815), doi:10.1029/2007GL030740.
- 5 Von Storch, H., and F.W. Zwiers, 1999: *Statistical Analysis in Climate Research*,
6 Cambridge University Press, 484 pp.
- 7 Wang, H., J.K.E. Schemm, A. Kumar, W. Wang, L. Long, M. Chelliah, G.D. Bell, and P.
8 Peng, 2009: A statistical forecast model for Atlantic seasonal hurricane activity
9 based on the NCEP dynamical seasonal forecast, *Journal of Climate*, **22**, 4481-
10 4500.
- 11 Yang, X. and co-authors (2012): A predictable AMO-like pattern in GFDL's fully-
12 coupled ensemble initialization and decadal forecasting system. *J. Climate*
13 (submitted).
- 14 Zhang, R., and T.L. Delworth, 2005: Simulated tropical response to a substantial
15 weakening of the Atlantic thermohaline circulation. *Journal of Climate*, **18**, 1853-
16 1860.
- 17 —, and —, 2006: Impact of Atlantic multidecadal oscillations on India/Sahel rainfall
18 and Atlantic hurricanes. *Geophysical Research Letters*, **33**, L17712,
19 doi:10.1029/2006GL026267.
- 20 —, and —, 2009: A new method for attributing climate variations over the Atlantic
21 hurricane basin's main development region. *Geophysical Research Letters*, **36**,
22 L06701, doi:10.1029/2009GL037260.

1 ———, and coauthors, 2012: Have aerosols caused the observed Atlantic Multidecadal
2 Variability? *Nature*, submitted.

3 Zhang, S., M.J. Harrison, A. Rosati, and A.T. Wittenberg, 2007: System design and
4 evaluation of coupled ensemble data assimilation for global oceanic climate
5 studies. *Mon. Wea. Rev.*, **135**, 3541–3564.

6 Zhang, S., and A. Rosati, 2010: An inflated ensemble filter for ocean data assimilation
7 with a biased coupled GCM. *Mon. Wea. Rev.*, **138**(10), 3905–3931.

8 Zhao, M., and I.M. Held, 2011: The response of tropical cyclone statistics to an increase
9 in CO₂ with fixed sea surface temperatures. *J. Climate*, **24**, 5353–5364.

10 ———, ———, S.-J. Lin, and G.A. Vecchi, 2009: Simulations of global hurricane
11 climatology, interannual variability, and response to global warming using a 50-
12 km resolution GCM. *J. Climate*, **22**, 6653–6678.

13 ———, ———, and G.A. Vecchi, 2010: Retrospective forecasts of the hurricane season using
14 a global atmospheric model assuming persistence of SST anomalies. *Mon. Wea.*
15 *Rev.*, **138**, 3858–3868.

16

Figure Captions:

Figure 1: Retrospective and future forecasts of hurricane frequency. Upper panels show the retrospective forecasts for five-year running hurricane frequency, lower panels focus on the nine-year running forecasts. Left panels show the results from uninitialized experiments, while the right panels show the results for initialized experiments. Black line shows the observed five-year (upper) and nine-year (lower) hurricane counts from the NOAA Hurricane Database (HURDAT; Jarvinen *et al.* 1984, MacAdie *et al.* 2009) that includes an adjustment for observing inhomogeneity prior to 1966 described in Vecchi and Knutson (2011). Retrospective forecasts are shown in: red line shows the forecasts from the GFDL-CM2.1 system, blue line shows the UKMO-DePreSys-PPE system, and the yellow line shows the two-system ensemble-mean.

Figure 2: Correlation for retrospective multi-year forecasts of North Atlantic hurricane frequency, with 90% uncertainty estimates. Each cluster of bars shows the retrospective correlation of multi-year hurricane frequency forecasts for Lead 2-6 years (left), Lead 6-10 years (middle) and Lead 2-10 years (right). Gray symbol is the correlation of the persistence of the five-year average count preceding the initialization of the model. Red symbols are for the GFDL-DecPre system, blue are for UKMO-DePreSys-PPE, and yellow is for the two system average. The initialized and uninitialized versions of each model are distinguished by different coloring. The sample correlation estimate is shown by the circle, the bars show the two-sided 90% uncertainty of a correlation given an underlying correlation with the value shown by the corresponding circle. Asterisk on top of a bar shows correlations that are significantly different from a null hypothesis of an

underlying correlation of zero at $p=0.1$, single-sided, with the effective degrees of freedom estimated as in Bretherton *et al.* (1999).

Figure 3: Mean Skill Score Squared (MSSS) of retrospective initialized multi-year hurricane frequency forecasts for various leads and models. Horizontal axis shows the MSSS against climatology, vertical axis shows the MSSS against the uninitialized forecasts; diagonal line indicates the one-to-one line. Left panel shows MSSS values for the five-year running-mean forecasts, right panel shows MSSS values for the nine-year running-mean forecasts. Circles show the values for the GFDL-DecPre system, squares for UKMO-DePreSys-PPE, and stars for the two-model ensemble mean. Different colors indicate different forecast leads.

Figure 4: Retrospective and future forecasts of the SST indices used for the hurricane emulator. Left panels show time-series of the five-year mean SSTA anomalies averaged over the global tropics (upper) and Atlantic hurricane main development region (lower), at lead 2-6. Black lines show observational estimates from HadISST.v1 (Rayner *et al.* 2003; solid) and ERSST.v3b (Smith *et al.* 2008; dotted). Colored lines show initialized (dashed) and uninitialized (solid) experiments from GFDL-DecPre (reds) and UKMO-DePreSys-PPE (blue). Right panels show the retrospective correlations of the forecasts at Lead 2-6 against the HadISST.v1 SST product.

Figure 5: Retrospective forecasts exploring the source of the initialized vs. uninitialized components. Left panel takes Atlantic MDR SST from initialized experiments and

tropical-mean SST from uninitialized, right panel takes tropical-mean SST from initialized experiments and Atlantic MDR SST from uninitialized experiments. The skill comes from the improvement of tropical Atlantic SST in the initialized experiments.

Figure 6: Retrospective forecasts arranged by verification and initialization date. Top panels (a and b) show the retrospective forecasts of five-year running hurricane averages for various leads, arranged so that each point on the time axis corresponds to the midpoint of the five-year interval over which the average is computed (*e.g.*, 1992 corresponds to the midpoint of the 1990-1994 average). Bottom panels (c and d) show the retrospective five-year forecasts for various leads arranged so that each point on the time axis corresponds to the date in which the model was initialized. Left panels are from the GFDL-CM2.1 forecasts, right panels are from the UKMO-DePreSys-PPE system. Dark line in the top panels shows the observed five-year running counts.

Figure 7: Empirical probability density function (PDF) estimates for the change in seasonal hurricane counts over the entire record and over the four years that preceded the 1994-1995 climate shift. The quantity explored is the difference in hurricane counts averaged over the five years following a given year with the counts of that year (*e.g.*, for 1991 it is the difference of hurricane counts averaged 1992-1996 with those in 1991); PDFs are estimated through Gaussian convolution with an *e*-folding scale of 2.5 hurricanes per year. Black lines are based on observations, blue lines on the forecasts with GFDL-DecPre, and red lines on the forecasts using UKMO-DePreSys; solid lines are computed over the 1961-2006 period, dashed lines over 1991-1994. The separation of

the solid and dashed black line is a reflection of the increase in storm counts that occurred in 1995. Notice that there is no tendency for forecasts initialized in the early-1990s to have indicate a tendency for intensification through the early years of the forecast: the forecast systems do not dynamically predict the occurrence of the 1994-1995 shift.

Figure 8: Retrospective forecasts of North Atlantic hurricane frequency after removing 1994-1995 shift in the mean from forecasts and verification (see Section III.A). Left panel shows the initialized forecasts at lead 2-6, right panels show the uninitialized experiments. Black line shows the observed counts, red line is from the GFDL-DecPre system, blue line is from UKMO-DePreSys-PPE and the yellow line is the two system average, all after removing the 1994-1995 shift in the mean.

Figure 9: Retrospective correlations of forecasts after removing 1994-1995 shift in the mean from forecasts and verification. Gray symbol is the correlation of the persistence of the five-year average count preceding the initialization of the model. Red symbols are for the GFDL-DecPre system, blue are for UKMO-DePreSys-PPE, and yellow is for the two system average. The initialized and uninitialized versions of each model are distinguished by different coloring. The sample correlation estimate is shown by the circle, the bars show the two-sided 90% uncertainty of a correlation given an underlying correlation with the value shown by the corresponding circle. Asterisk on top of a bar shows correlations that are significantly different from a null hypothesis of an underlying correlation of zero at $p=0.1$, single-sided, with the effective degrees of freedom estimated as in Bretherton *et al.* (1999).

1
2

Forecast system	Underlying GCM	Initialization Procedure	Ensemble Type	Initialization Date	Treatment of Volcanoes
GFDL-CM2.1 DecPre (Rosati <i>et al.</i> 2012; Yang <i>et al.</i> 2012)	GFDL-CM2.1 (Delworth <i>et al.</i> 2006)	Fully Coupled Ensemble Kalman Filter (Zhang <i>et al.</i> 2007), full variable assimilation	Ten ensemble members from the EnKF assimilation	1-January of each year 1960-2011.	Future volcanoes included in radiative forcing
UKMO DepPreSys-PPE (Smith <i>et al.</i> 2007; Smith <i>et al.</i> 2010)	HadCM3 (Gordon <i>et al.</i> 2000)	Atmospheric and oceanic conditions relaxed to observations. Ocean anomaly initialization. (Smith and Murphy 2007)	Nine ensemble member perturbed physics ensemble (PPE)	1 November of each year 1960-2005.	Forcing from past volcanic forcing included

3 **Table 1:** Summary of the two dynamical multi-year experimental forecast systems
4 explored in this manuscript.

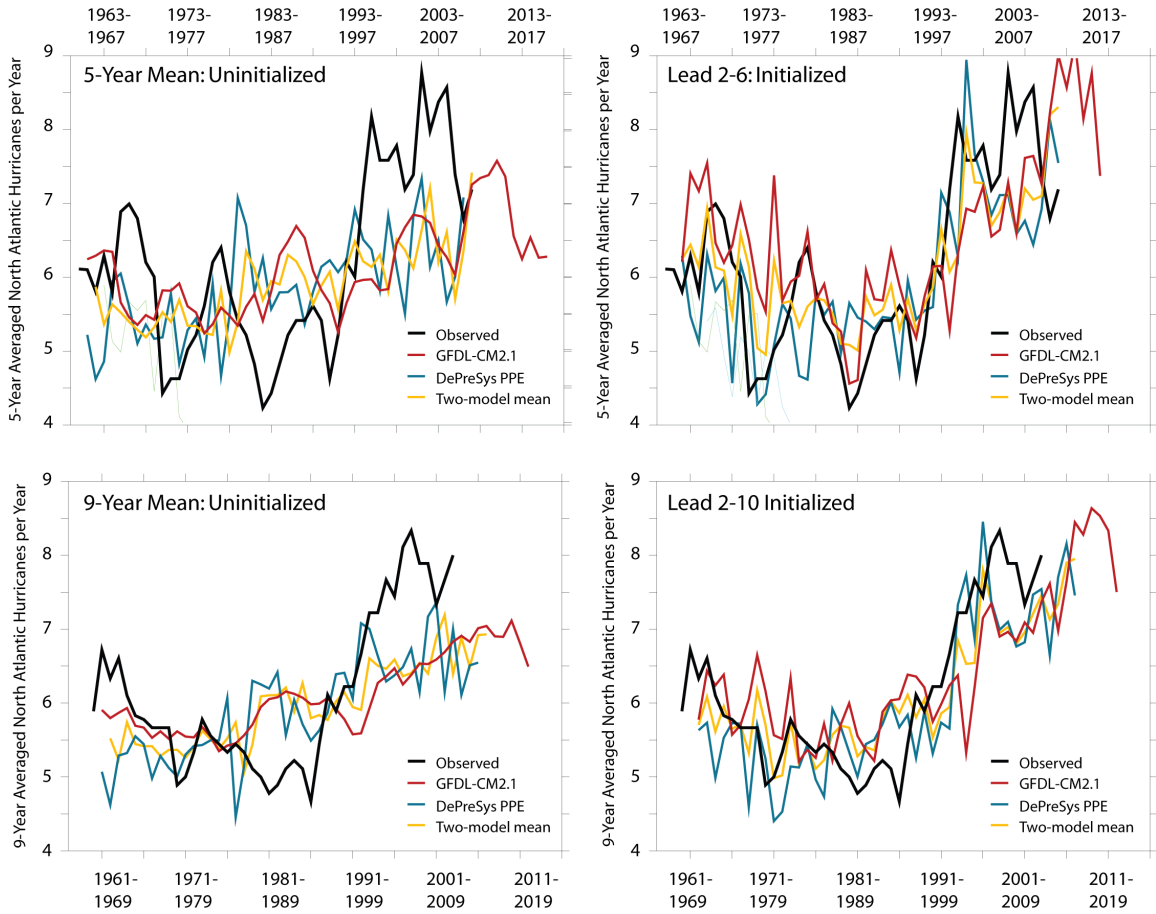


Figure 1: Retrospective and future forecasts of hurricane frequency. Upper panels show the retrospective forecasts for five-year running hurricane frequency, lower panels focus on the nine-year running forecasts. Left panels show the results from uninitialized experiments, while the right panels show the results for initialized experiments. Black line shows the observed five-year hurricane counts from the NOAA Hurricane Database (HURDAT; Jarvinen *et al.* 1984, MacAdie *et al.* 2009) that includes an adjustment for observing inhomogeneity prior to 1966 described in Vecchi and Knutson (2011). Retrospective forecasts are shown in: red line shows the forecasts from the GFDL-CM2.1 system, blue line shows the UKMO-DePreSys-PPE system, and the yellow line shows the two-system ensemble-mean.

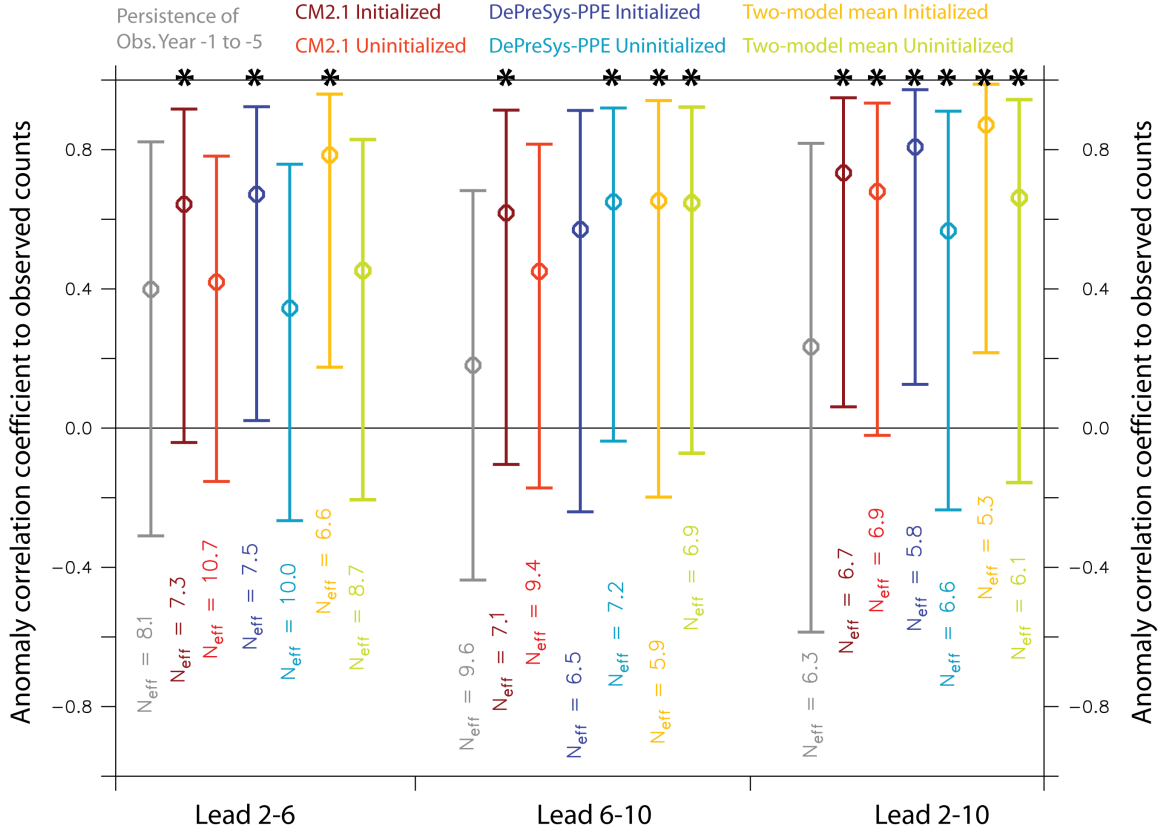


Figure 2: Correlation for retrospective multi-year forecasts of North Atlantic hurricane frequency, with 90% uncertainty estimates. Each cluster of bars shows the retrospective correlation of multi-year hurricane frequency forecasts for lead 2-6 years (left), lead 6-10 years (middle) and lead 2-10 years (right). Gray symbol is the correlation of the persistence of the five-year average count preceding the initialization of the model. Red symbols are for the GFDL-DecPre system, blue are for UKMO-DePreSys-PPE, and yellow is for the two system average. The initialized and uninitialized versions of each model are distinguished by different coloring. The sample correlation estimate is shown by the circle, the bars show the two-sided 90% uncertainty of a correlation given an underlying correlation with the value shown by the corresponding circle. Asterisk on top of a bar shows correlations that are significantly different from a null hypothesis of an underlying correlation of zero at $p=0.1$, single-sided, with the effective degrees of freedom estimated as in Bretherton *et al.* (1999).

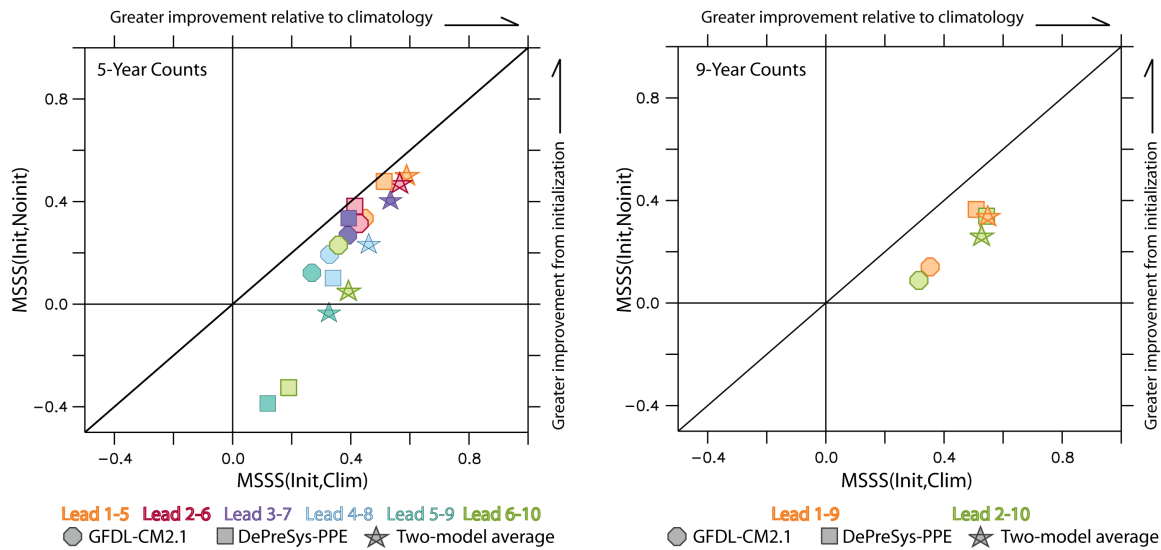


Figure 3: Mean Skill Score Squared (MSSS) of retrospective initialized multi-year hurricane frequency forecasts for various leads and models. Horizontal axis shows the MSSS against climatology, vertical axis shows the MSSS against the uninitialized forecasts; diagonal line indicates the one-to-one line. Left panel shows MSSS values for the five-year running-mean forecasts, right panel shows MSSS values for the nine-year running-mean forecasts. Circles show the values for the GFDL-DecPre system, squares for UKMO-DecPreSys-PPE, and stars for the two-model ensemble mean. Different colors indicate different forecast leads.

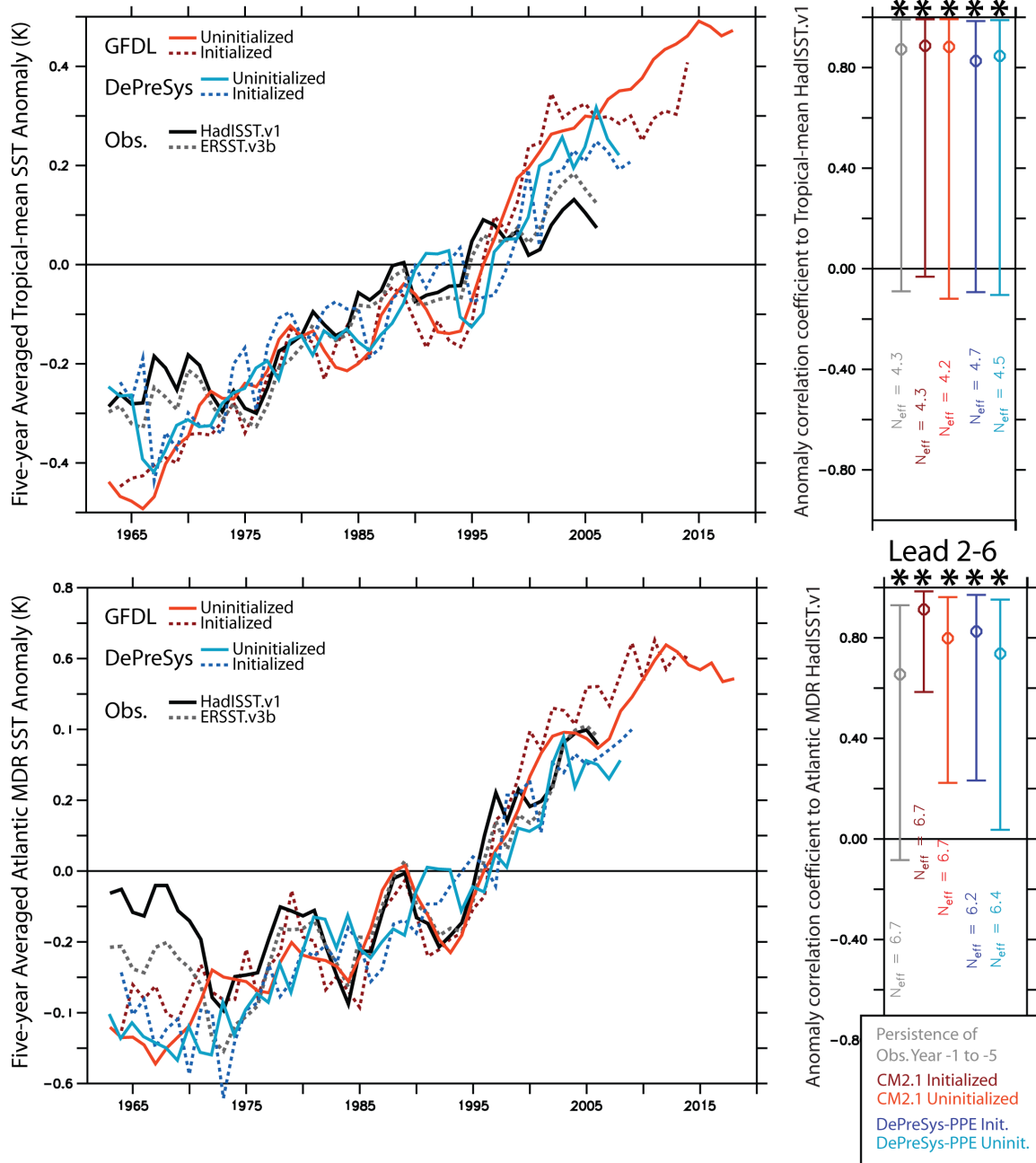
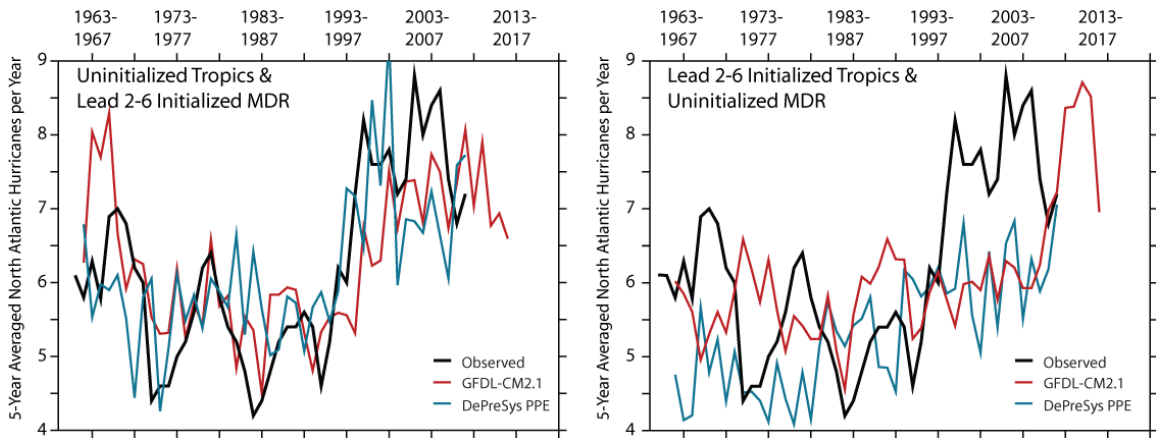


Figure 4: Retrospective and future forecasts of the SST indices used for the hurricane emulator. Left panels show time-series of the five-year mean SST anomalies averaged over the global tropics (upper) and Atlantic hurricane main development region (lower), at lead 2-6. Black lines show observational estimates from HadISST.v1 (Rayner *et al.* 2003; solid) and ERSST.v3b (Smith *et al.* 2008; dotted). Colored lines show initialized (dashed) and uninitialized (solid) experiments from GFDL-DecPre (reds) and UKMO-DePreSys-PPE (blue). Right panels show the retrospective correlations of the forecasts at lead 2-6 against the HadISST.v1 SST product.

1

2



3

4

5

6

7

8

Figure 5: Retrospective forecasts exploring the source of the initialized vs. uninitialized components. Left panel takes Atlantic MDR SST from initialized experiments and tropical-mean SST from uninitialized experiments, right panel takes tropical-mean SST from initialized experiments and Atlantic MDR SST from uninitialized experiments. The skill comes from the improvement of tropical Atlantic SST in the initialized experiments.

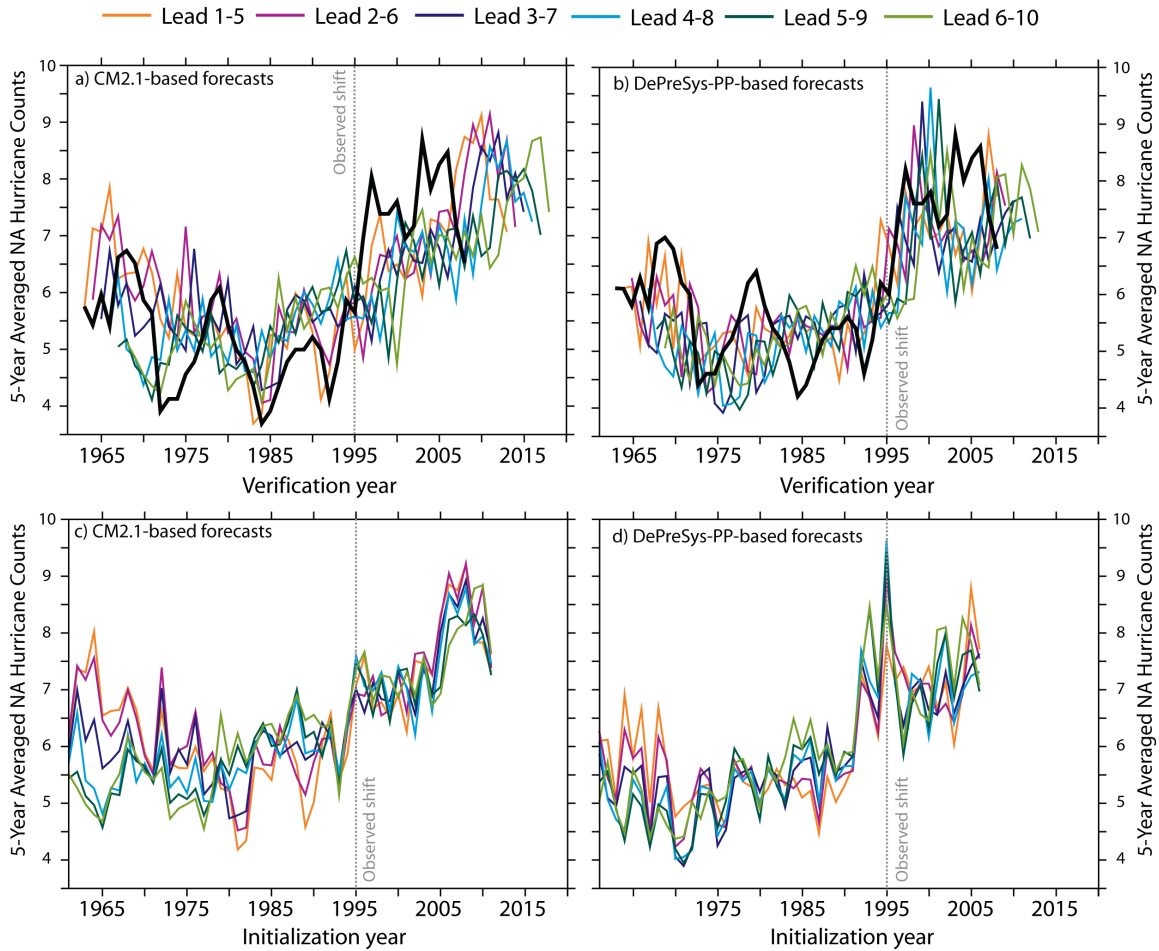
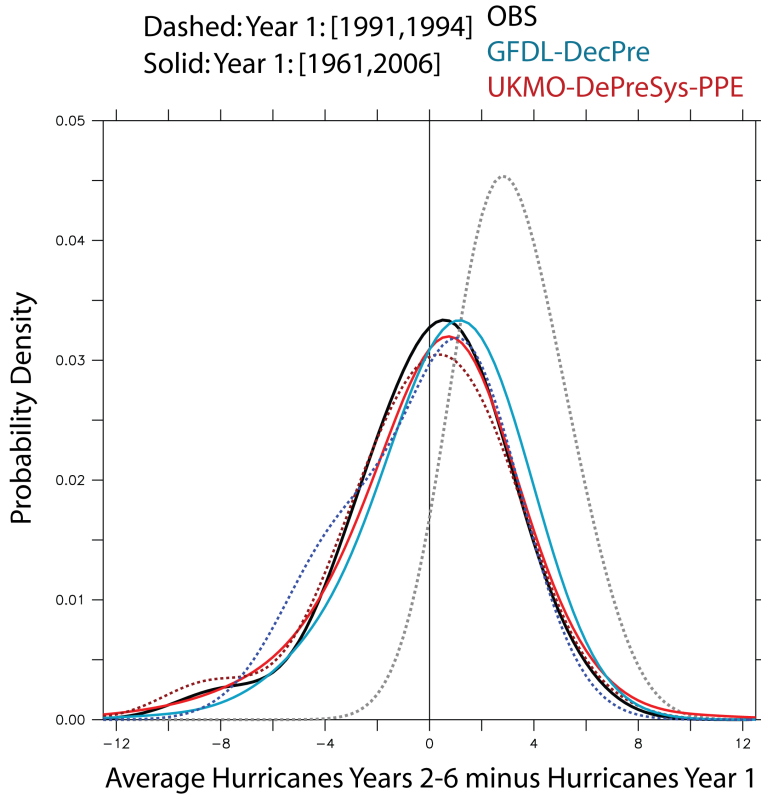


Figure 6: Retrospective forecasts arranged by verification and initialization date. Top panels (a and b) show the retrospective forecasts of five-year running hurricane averages for various leads, arranged so that each point on the time axis corresponds to the midpoint of the five-year interval over which the average is computed (*e.g.*, 1992 corresponds to the midpoint of the 1990-1994 average). Bottom panels (c and d) show the retrospective five-year forecasts for various leads arranged so that each point on the time axis corresponds to the date in which the model was initialized. Left panels are from the GFDL-CM2.1 forecasts, right panels are from the UKMO-DePreSys-PPE system. Dark line in the top panels shows the observed five-year running counts.



1 Average Hurricanes Years 2-6 minus Hurricanes Year 1
2 **Figure 7:** Empirical probability density function (PDF) estimates for the change in
3 seasonal hurricane counts over the entire record and over the four years that preceded the
4 1994-1995 climate shift. The quantity explored is the difference in hurricane counts
5 averaged over the five years following a given year with the counts of that year (*e.g.*, for
6 1991 it is the difference of hurricane counts averaged 1992-1996 with those in 1991);
7 PDFs are estimated through Gaussian convolution with an e -folding scale of 2.5
8 hurricanes per year. Black lines are based on observations, blue lines on the forecasts
9 with GFDL-DecPre, and red lines on the forecasts using UKMO-DePreSys; solid lines
10 are computed over the 1961-2006 period, dashed lines over 1991-1994. PDFs of the
11 models are based on the various ensemble members. The separation of the solid and
12 dashed black lines is a reflection of the increase in storm counts that occurred in 1995.
13 Notice that there is no tendency for forecasts initialized in the early-1990s to have
14 indicate a tendency for frequency increase through the early years of the forecast: the
15 forecast systems do not dynamically predict the occurrence of the 1994-1995 shift.

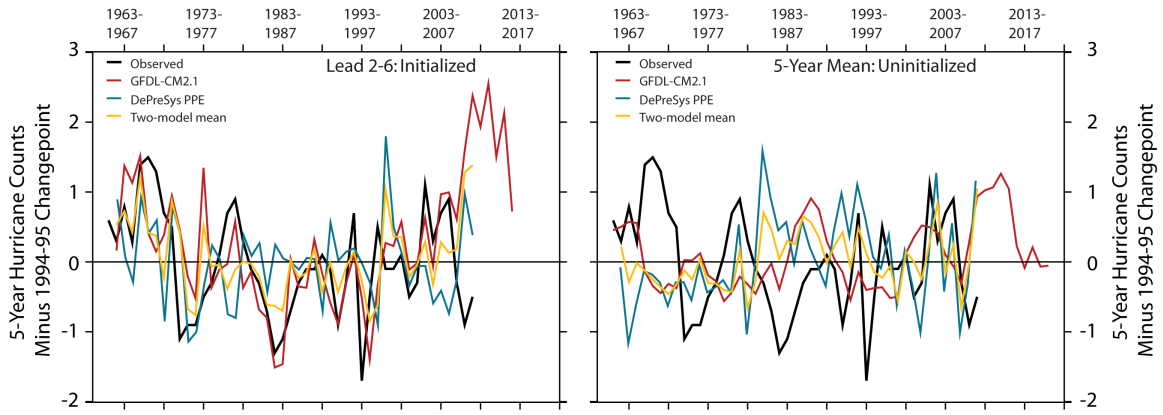


Figure 8: Retrospective forecasts of North Atlantic hurricane frequency after removing 1994-1995 shift in the mean from forecasts and verification (see Section III.A). Left panel shows the initialized forecasts at lead 2-6, right panel shows the uninitialized experiments. Black line shows the observed counts, red line is from the GFDL-DecPre system, blue line is from UKMO-DePreSys-PPE and the yellow line is the two system average, all after removing the 1994-1995 shift in the mean.

Persistence of Obs. Year -1 to -5 CM2.1 Initialized DePreSys-PPE Initialized Two-model mean Initialized
 CM2.1 Uninitialized DePreSys-PPE Uninitialized Two-model mean Uninitialized

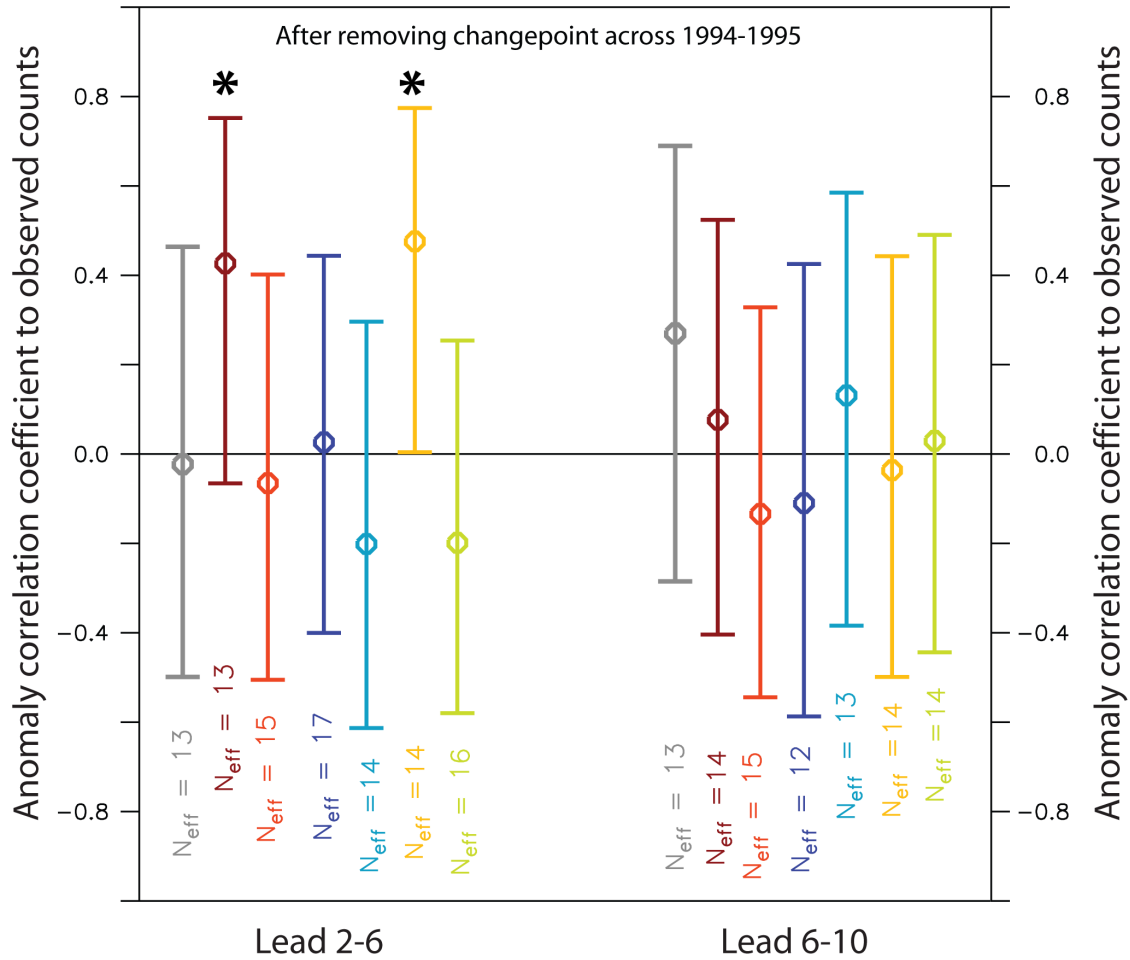


Figure 9: Retrospective correlations of forecasts after removing 1994-1995 shift in the mean from forecasts and verification. Gray symbol is the correlation of the persistence of the five-year average count preceding the initialization of the model. Red symbols are for the GFDL-DecPre system, blue are for UKMO-DePreSys-PPE, and yellow is for the two system average. The initialized and uninitialized versions of each model are distinguished by different coloring. The sample correlation estimate is shown by the circle, the bars show the two-sided 90% uncertainty of a correlation given an underlying correlation with the value shown by the corresponding circle. Asterisk on top of a bar shows correlations that are significantly different from a null hypothesis of an underlying correlation of zero at $p=0.1$, single-sided, with the effective degrees of freedom estimated as in Bretherton *et al.* (1999).